

CS224n: Natural Language Processing with Deep Learning¹

Lecture Notes: Part I

Word Vectors I: Introduction, SVD and Word2Vec²

Winter 2019

¹ Course Instructors: Christopher Manning, Richard Socher

² Authors: Francois Chaubard, Michael Fang, Guillaume Genthial, Rohit Mundra, Richard Socher

Keyphrases: Natural Language Processing. Word Vectors. Singular Value Decomposition. Skip-gram. Continuous Bag of Words (CBOW). Negative Sampling. Hierarchical Softmax. Word2Vec.

This set of notes begins by introducing the concept of Natural Language Processing (NLP) and the problems NLP faces today. We then move forward to discuss the concept of representing words as numeric vectors. Lastly, we discuss popular approaches to designing word vectors.

1 Introduction to Natural Language Processing

We begin with a general discussion of what is NLP.

1.1 What is so special about NLP?

What's so special about human (natural) language? Human language is a system specifically constructed to convey meaning, and is not produced by a physical manifestation of any kind. In that way, it is very different from vision or any other machine learning task.

Most words are just symbols for an extra-linguistic entity : the word is a *signifier* that maps to a *signified* (idea or thing).

For instance, the word "rocket" refers to the concept of a rocket, and by extension can designate an instance of a rocket. There are some exceptions, when we use words and letters for expressive signaling, like in "Whooompaa". On top of this, the symbols of language can be encoded in several modalities : voice, gesture, writing, etc that are transmitted via *continuous* signals to the brain, which itself appears to encode things in a continuous manner. (A lot of work in philosophy of language and linguistics has been done to conceptualize human language and distinguish words from their references, meanings, etc. Among others, see works by Wittgenstein, Frege, Russell and Mill.)

Natural language is a discrete/symbolic/categorical system

1.2 Examples of tasks

There are different levels of tasks in NLP, from speech processing to semantic interpretation and discourse processing. The goal of NLP is to be able to design algorithms to allow computers to "understand"

natural language in order to perform some task. Example tasks come in varying level of difficulty:

Easy

- Spell Checking
- Keyword Search
- Finding Synonyms

Medium

- Parsing information from websites, documents, etc.

Hard

- Machine Translation (e.g. Translate Chinese text to English)
- Semantic Analysis (What is the meaning of query statement?)
- Coreference (e.g. What does "he" or "it" refer to given a document?)
- Question Answering (e.g. Answering Jeopardy questions).

1.3 How to represent words?

The first and arguably most important common denominator across all NLP tasks is how we represent words as input to any of our models. Much of the earlier NLP work that we will not cover treats words as atomic symbols. To perform well on most NLP tasks we first need to have some notion of similarity and difference between words. With word vectors, we can quite easily encode this ability in the vectors themselves (using distance measures such as Jaccard, Cosine, Euclidean, etc).

2 Word Vectors

There are an estimated 13 million tokens for the English language but are they all completely unrelated? Feline to cat, hotel to motel? I think not. Thus, we want to encode word tokens each into some vector that represents a point in some sort of "word" space. This is paramount for a number of reasons but the most intuitive reason is that perhaps there actually exists some N -dimensional space (such that $N \ll 13$ million) that is sufficient to encode all semantics of our language. Each dimension would encode some meaning that we transfer using speech. For instance, semantic dimensions might

indicate tense (past vs. present vs. future), count (singular vs. plural), and gender (masculine vs. feminine).

So let's dive into our first word vector and arguably the most simple, the **one-hot vector**: Represent every word as an $\mathbb{R}^{|V| \times 1}$ vector with all 0s and one 1 at the index of that word in the sorted english language. In this notation, $|V|$ is the size of our vocabulary. Word vectors in this type of encoding would appear as the following:

$$w^{aardvark} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, w^a = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, w^{at} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots w^{zebra} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

We represent each word as a completely independent entity. As we previously discussed, this word representation does not give us directly any notion of similarity. For instance,

$$(w^{hotel})^T w^{motel} = (w^{hotel})^T w^{cat} = 0$$

So maybe we can try to reduce the size of this space from $\mathbb{R}^{|V|}$ to something smaller and thus find a subspace that encodes the relationships between words.

One-hot vector: Represent every word as an $\mathbb{R}^{|V| \times 1}$ vector with all 0s and one 1 at the index of that word in the sorted english language.

Fun fact: The term "one-hot" comes from digital circuit design, meaning "a group of bits among which the legal combinations of values are only those with a single high (1) bit and all the others low (0)".

Denotational semantics: The concept of representing an idea as a symbol (a word or a one-hot vector). It is sparse and cannot capture similarity. This is a "localist" representation.

3 SVD Based Methods

For this class of methods to find word embeddings (otherwise known as word vectors), we first loop over a massive dataset and accumulate word co-occurrence counts in some form of a matrix X , and then perform Singular Value Decomposition on X to get a USV^T decomposition. We then use the rows of U as the word embeddings for all words in our dictionary. Let us discuss a few choices of X .

3.1 Word-Document Matrix

As our first attempt, we make the bold conjecture that words that are related will often appear in the same documents. For instance, "banks", "bonds", "stocks", "money", etc. are probably likely to appear together. But "banks", "octopus", "banana", and "hockey" would probably not consistently appear together. We use this fact to build a word-document matrix, X in the following manner: Loop over billions of documents and for each time word i appears in document j , we add one to entry X_{ij} . This is obviously a very large matrix ($\mathbb{R}^{|V| \times M}$) and it scales with the number of documents (M). So perhaps we can try something better.

Distributional semantics: The concept of representing the meaning of a word based on the context in which it usually appears. It is dense and can better capture similarity.

3.2 Window based Co-occurrence Matrix

The same kind of logic applies here however, the matrix X stores co-occurrences of words thereby becoming an affinity matrix. In this method we count the number of times each word appears inside a window of a particular size around the word of interest. We calculate this count for all the words in corpus. We display an example below. Let our corpus contain just three sentences and the window size be 1:

1. I enjoy flying.
2. I like NLP.
3. I like deep learning.

The resulting counts matrix will then be:

$$X = \begin{matrix} & \begin{matrix} I & like & enjoy & deep & learning & NLP & flying & . \end{matrix} \\ \begin{matrix} I \\ like \\ enjoy \\ deep \\ learning \\ NLP \\ flying \\ . \end{matrix} & \begin{bmatrix} 0 & 2 & 1 & 0 & 0 & 0 & 0 & 0 \\ 2 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \end{bmatrix} \end{matrix}$$

Using Word-Word Co-occurrence Matrix:

- Generate $|V| \times |V|$ co-occurrence matrix, X .
- Apply SVD on X to get $X = USV^T$.
- Select the first k columns of U to get a k -dimensional word vectors.
- $\frac{\sum_{i=1}^k \sigma_i}{\sum_{i=1}^{|V|} \sigma_i}$ indicates the amount of variance captured by the first k dimensions.

3.3 Applying SVD to the cooccurrence matrix

We now perform SVD on X , observe the singular values (the diagonal entries in the resulting S matrix), and cut them off at some index k based on the desired percentage variance captured:

$$\frac{\sum_{i=1}^k \sigma_i}{\sum_{i=1}^{|V|} \sigma_i}$$

We then take the submatrix of $U_{1:|V|,1:k}$ to be our word embedding matrix. This would thus give us a k -dimensional representation of every word in the vocabulary.

Applying SVD to X :

$$\begin{matrix} |V| \\ \left[\begin{array}{c} X \end{array} \right] \end{matrix} = \begin{matrix} |V| \\ \left[\begin{array}{c|c|c} u_1 & u_2 & \dots \end{array} \right] \end{matrix} \begin{matrix} |V| \\ \left[\begin{array}{c} \sigma_1 & 0 & \dots \\ 0 & \sigma_2 & \dots \\ \vdots & \vdots & \ddots \end{array} \right] \end{matrix} \begin{matrix} |V| \\ \left[\begin{array}{c} - & v_1 & - \\ - & v_2 & - \\ \vdots & \vdots & \vdots \end{array} \right] \end{matrix}$$

Reducing dimensionality by selecting first k singular vectors:

$$|V| \begin{bmatrix} | & & \\ \hat{X} & & \\ | & & \end{bmatrix} = |V| \begin{bmatrix} | & & \\ u_1 & u_2 & \dots \\ | & & \end{bmatrix}^k \begin{bmatrix} \sigma_1 & 0 & \dots \\ 0 & \sigma_2 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}^k \begin{bmatrix} - & v_1 & - \\ - & v_2 & - \\ \vdots & \vdots & \end{bmatrix}$$

Both of these methods give us word vectors that are more than sufficient to encode semantic and syntactic (part of speech) information but are associated with many other problems:

- The dimensions of the matrix change very often (new words are added very frequently and corpus changes in size).
- The matrix is extremely sparse since most words do not co-occur.
- The matrix is very high dimensional in general ($\approx 10^6 \times 10^6$)
- Quadratic cost to train (i.e. to perform SVD)
- Requires the incorporation of some hacks on X to account for the drastic imbalance in word frequency

SVD based methods do not scale well for big matrices and it is hard to incorporate new words or documents. Computational cost for a $m \times n$ matrix is $O(mn^2)$

Some solutions exist to resolve some of the issues discussed above:

However, count-based methods make an efficient use of the statistics

- Ignore function words such as "the", "he", "has", etc.
- Apply a ramp window – i.e. weight the co-occurrence count based on distance between the words in the document.
- Use Pearson correlation and set negative counts to 0 instead of using just raw count.

As we see in the next section, iteration based methods solve many of these issues in a far more elegant manner.

4 Iteration Based Methods - Word2vec

Let us step back and try a new approach. Instead of computing and storing global information about some huge dataset (which might be billions of sentences), we can try to create a model that will be able to learn one iteration at a time and eventually be able to encode the probability of a word given its context.

For an overview of Word2vec, a note map can be found here : <https://myndbook.com/view/4900>

A detailed summary of word2vec models can also be found here [Rong, 2014]

The idea is to design a model whose parameters are the word vectors. Then, train the model on a certain objective. At every iteration we run our model, evaluate the errors, and follow an update rule that has some notion of penalizing the model parameters that caused the error. Thus, we learn our word vectors. This idea is a very old

Iteration-based methods capture co-occurrence of words one at a time instead of capturing all cooccurrence counts directly like in SVD methods.

one dating back to 1986. We call this method "backpropagating" the errors (see [Rumelhart et al., 1988]). The simpler the model and the task, the faster it will be to train it.

Several approaches have been tested. [Collobert et al., 2011] design models for NLP whose first step is to transform each word in a vector. For each special task (Named Entity Recognition, Part-of-Speech tagging, etc.) they train not only the model's parameters but also the vectors and achieve great performance, while computing good word vectors! Other interesting reading would be [Bengio et al., 2003].

In this class, we will present a simpler, more recent, probabilistic method by [Mikolov et al., 2013] : word2vec. Word2vec is a software package that actually includes :

- **2 algorithms**: continuous bag-of-words (CBOW) and skip-gram. CBOW aims to predict a center word from the surrounding context in terms of word vectors. Skip-gram does the opposite, and predicts the *distribution* (probability) of context words from a center word.

- **2 training methods**: negative sampling and hierarchical softmax. Negative sampling defines an objective by sampling *negative* examples, while hierarchical softmax defines an objective using an efficient tree structure to compute probabilities for all the vocabulary.

4.1 Language Models (Unigrams, Bigrams, etc.)

First, we need to create such a model that will assign a probability to a sequence of tokens. Let us start with an example:

"The cat jumped over the puddle."

A good language model will give this sentence a high probability because this is a completely valid sentence, syntactically and semantically. Similarly, the sentence "stock boil fish is toy" should have a very low probability because it makes no sense. Mathematically, we can call this probability on any given sequence of n words:

$$P(w_1, w_2, \dots, w_n)$$

We can take the unary language model approach and break apart this probability by assuming the word occurrences are completely independent:

$$P(w_1, w_2, \dots, w_n) = \prod_{i=1}^n P(w_i)$$

However, we know this is a bit ludicrous because we know the next word is highly contingent upon the previous sequence of words. And the silly sentence example might actually score highly. So perhaps we let the probability of the sequence depend on the pairwise

Context of a word:

The context of a word is the set of m surrounding words. For instance, the $m = 2$ context of the word "fox" in the sentence "The quick brown fox jumped over the lazy dog" is {"quick", "brown", "jumped", "over"}.

This model relies on a very important hypothesis in linguistics, *distributional similarity*, the idea that similar words have similar context.

Unigram model:

$$P(w_1, w_2, \dots, w_n) = \prod_{i=1}^n P(w_i)$$

probability of a word in the sequence and the word next to it. We call this the bigram model and represent it as:

$$P(w_1, w_2, \dots, w_n) = \prod_{i=2}^n P(w_i | w_{i-1})$$

Again this is certainly a bit naive since we are only concerning ourselves with pairs of neighboring words rather than evaluating a whole sentence, but as we will see, this representation gets us pretty far along. Note in the Word-Word Matrix with a context of size 1, we basically can learn these pairwise probabilities. But again, this would require computing and storing global information about a massive dataset.

Now that we understand how we can think about a sequence of tokens having a probability, let us observe some example models that could learn these probabilities.

4.2 Continuous Bag of Words Model (CBOW)

One approach is to treat {"The", "cat", "over", "the", "puddle"} as a context and from these words, be able to predict or generate the center word "jumped". This type of model we call a Continuous Bag of Words (CBOW) Model.

Let's discuss the CBOW Model above in greater detail. First, we set up our known parameters. Let the known parameters in our model be the sentence represented by one-hot word vectors. The input one hot vectors or context we will represent with an $x^{(c)}$. And the output as $y^{(c)}$ and in the CBOW model, since we only have one output, so we just call this y which is the one hot vector of the known center word. Now let's define our unknowns in our model.

We create two matrices, $\mathcal{V} \in \mathbb{R}^{n \times |V|}$ and $\mathcal{U} \in \mathbb{R}^{|V| \times n}$. Where n is an arbitrary size which defines the size of our embedding space. \mathcal{V} is the input word matrix such that the i -th column of \mathcal{V} is the n -dimensional embedded vector for word w_i when it is an input to this model. We denote this $n \times 1$ vector as v_i . Similarly, \mathcal{U} is the output word matrix. The j -th row of \mathcal{U} is an n -dimensional embedded vector for word w_j when it is an output of the model. We denote this row of \mathcal{U} as u_j . Note that we do in fact learn two vectors for every word w_i (i.e. input word vector v_i and output word vector u_i).

Bigram model:

$$P(w_1, w_2, \dots, w_n) = \prod_{i=2}^n P(w_i | w_{i-1})$$

CBOW Model:

Predicting a center word from the surrounding context

For each word, we want to learn 2 vectors

- v : (input vector) when the word is in the context
- u : (output vector) when the word is in the center

Notation for CBOW Model:

- w_i : Word i from vocabulary V
- $\mathcal{V} \in \mathbb{R}^{n \times |V|}$: Input word matrix
- v_i : i -th column of \mathcal{V} , the input vector representation of word w_i
- $\mathcal{U} \in \mathbb{R}^{|V| \times n}$: Output word matrix
- u_i : i -th row of \mathcal{U} , the output vector representation of word w_i

We breakdown the way this model works in these steps:

1. We generate our one hot word vectors for the input context of size m : $(x^{(c-m)}, \dots, x^{(c-1)}, x^{(c+1)}, \dots, x^{(c+m)}) \in \mathbb{R}^{|V|}$.

2. We get our embedded word vectors for the context ($v_{c-m} = \mathcal{V}x^{(c-m)}, v_{c-m+1} = \mathcal{V}x^{(c-m+1)}, \dots, v_{c+m} = \mathcal{V}x^{(c+m)} \in \mathbb{R}^n$)
3. Average these vectors to get $\hat{v} = \frac{v_{c-m} + v_{c-m+1} + \dots + v_{c+m}}{2m} \in \mathbb{R}^n$
4. Generate a score vector $z = \mathcal{U}\hat{v} \in \mathbb{R}^{|V|}$. As the dot product of similar vectors is higher, it will push similar words close to each other in order to achieve a high score.
5. Turn the scores into probabilities $\hat{y} = \text{softmax}(z) \in \mathbb{R}^{|V|}$.
6. We desire our probabilities generated, $\hat{y} \in \mathbb{R}^{|V|}$, to match the true probabilities, $y \in \mathbb{R}^{|V|}$, which also happens to be the one hot vector of the actual word.

So now that we have an understanding of how our model would work if we had a \mathcal{V} and \mathcal{U} , how would we learn these two matrices? Well, we need to create an objective function. Very often when we are trying to learn a probability from some true probability, we look to information theory to give us a measure of the distance between two distributions. Here, we use a popular choice of distance/loss measure, cross entropy $H(\hat{y}, y)$.

The intuition for the use of cross-entropy in the discrete case can be derived from the formulation of the loss function:

$$H(\hat{y}, y) = - \sum_{j=1}^{|V|} y_j \log(\hat{y}_j)$$

Let us concern ourselves with the case at hand, which is that y is a one-hot vector. Thus we know that the above loss simplifies to simply:

$$H(\hat{y}, y) = -y_i \log(\hat{y}_i)$$

In this formulation, c is the index where the correct word's one hot vector is 1. We can now consider the case where our prediction was perfect and thus $\hat{y}_c = 1$. We can then calculate $H(\hat{y}, y) = -1 \log(1) = 0$. Thus, for a perfect prediction, we face no penalty or loss. Now let us consider the opposite case where our prediction was very bad and thus $\hat{y}_c = 0.01$. As before, we can calculate our loss to be $H(\hat{y}, y) = -1 \log(0.01) \approx 4.605$. We can thus see that for probability distributions, cross entropy provides us with a good measure of distance. We thus formulate our optimization objective as:

The **softmax** is an operator that we'll use very frequently. It transforms a vector into a vector whose i -th component is $\frac{e^{y_i}}{\sum_{k=1}^{|V|} e^{y_k}}$.

- exponentiate to make positive
- Dividing by $\sum_{k=1}^{|V|} e^{y_k}$ normalizes the vector ($\sum_{k=1}^n \hat{y}_k = 1$) to give probability

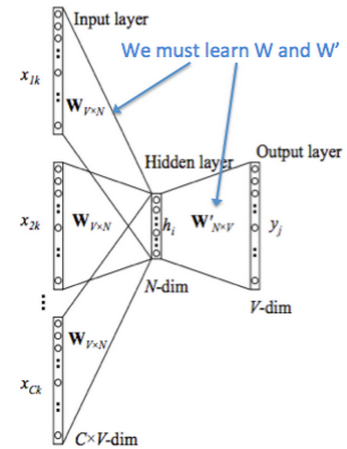


Figure 1: This image demonstrates how CBOW works and how we must learn the transfer matrices

$\hat{y} \mapsto H(\hat{y}, y)$ is minimum when $\hat{y} = y$. Then, if we found a \hat{y} such that $H(\hat{y}, y)$ is close to the minimum, we have $\hat{y} \approx y$. This means that our model is very good at predicting the center word!

To learn the vectors (the matrices U and V) CBOW defines a cost that measures how good it is at predicting the center word. Then, we optimize this cost by updating the matrices U and V thanks to stochastic gradient descent

$$\begin{aligned}
 \text{minimize } J &= -\log P(w_c | w_{c-m}, \dots, w_{c-1}, w_{c+1}, \dots, w_{c+m}) \\
 &= -\log P(u_c | \hat{v}) \\
 &= -\log \frac{\exp(u_c^T \hat{v})}{\sum_{j=1}^{|V|} \exp(u_j^T \hat{v})} \\
 &= -u_c^T \hat{v} + \log \sum_{j=1}^{|V|} \exp(u_j^T \hat{v})
 \end{aligned}$$

We use stochastic gradient descent to update all relevant word vectors u_c and v_j .

4.3 Skip-Gram Model

Another approach is to create a model such that given the center word "jumped", the model will be able to predict or generate the surrounding words "The", "cat", "over", "the", "puddle". Here we call the word "jumped" the context. We call this type of model a Skip-Gram model.

Let's discuss the Skip-Gram model above. The setup is largely the same but we essentially swap our x and y i.e. x in the CBOW are now y and vice-versa. The input one hot vector (center word) we will represent with an x (since there is only one). And the output vectors as $y^{(j)}$. We define \mathcal{V} and \mathcal{U} the same as in CBOW.

We breakdown the way this model works in these 6 steps:

1. We generate our one hot input vector $x \in \mathbb{R}^{|V|}$ of the center word.
2. We get our embedded word vector for the center word $v_c = \mathcal{V}x \in \mathbb{R}^n$
3. Generate a score vector $z = \mathcal{U}v_c$.
4. Turn the score vector into probabilities, $\hat{y} = \text{softmax}(z)$. Note that $\hat{y}_{c-m}, \dots, \hat{y}_{c-1}, \hat{y}_{c+1}, \dots, \hat{y}_{c+m}$ are the probabilities of observing each context word.
5. We desire our probability vector generated to match the true probabilities which is $y^{(c-m)}, \dots, y^{(c-1)}, y^{(c+1)}, \dots, y^{(c+m)}$, the one hot vectors of the actual output.

As in CBOW, we need to generate an objective function for us to evaluate the model. A key difference here is that we invoke a Naive Bayes assumption to break out the probabilities. If you have not seen this before, then simply put, it is a strong (naive) conditional

Stochastic gradient descent (SGD) computes gradients for a window and updates the parameters

$$\mathcal{U}_{new} \leftarrow \mathcal{U}_{old} - \alpha \nabla_{\mathcal{U}} J$$

$$\mathcal{V}_{old} \leftarrow \mathcal{V}_{old} - \alpha \nabla_{\mathcal{V}} J$$

Skip-Gram Model:

Predicting surrounding context words given a center word

Notation for Skip-Gram Model:

- w_i : Word i from vocabulary V
- $\mathcal{V} \in \mathbb{R}^{n \times |V|}$: Input word matrix
- v_i : i -th column of \mathcal{V} , the input vector representation of word w_i
- $\mathcal{U} \in \mathbb{R}^{n \times |V|}$: Output word matrix
- u_i : i -th row of \mathcal{U} , the output vector representation of word w_i

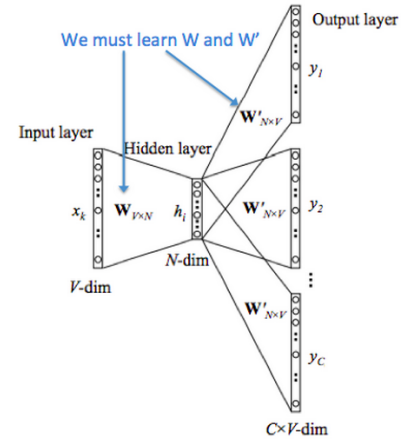


Figure 2: This image demonstrates how Skip-Gram works and how we must learn the transfer matrices

independence assumption. In other words, given the center word, all output words are completely independent.

$$\begin{aligned}
 \text{minimize } J &= -\log P(w_{c-m}, \dots, w_{c-1}, w_{c+1}, \dots, w_{c+m} | w_c) \\
 &= -\log \prod_{j=0, j \neq m}^{2m} P(w_{c-m+j} | w_c) \\
 &= -\log \prod_{j=0, j \neq m}^{2m} P(u_{c-m+j} | v_c) \\
 &= -\log \prod_{j=0, j \neq m}^{2m} \frac{\exp(u_{c-m+j}^T v_c)}{\sum_{k=1}^{|V|} \exp(u_k^T v_c)} \\
 &= -\sum_{j=0, j \neq m}^{2m} u_{c-m+j}^T v_c + 2m \log \sum_{k=1}^{|V|} \exp(u_k^T v_c)
 \end{aligned}$$

With this objective function, we can compute the gradients with respect to the unknown parameters and at each iteration update them via Stochastic Gradient Descent.

Note that

$$\begin{aligned}
 J &= -\sum_{j=0, j \neq m}^{2m} \log P(u_{c-m+j} | v_c) \\
 &= \sum_{j=0, j \neq m}^{2m} H(\hat{y}, y_{c-m+j})
 \end{aligned}$$

where $H(\hat{y}, y_{c-m+j})$ is the cross-entropy between the probability vector \hat{y} and the one-hot vector y_{c-m+j} .

Only one probability vector \hat{y} is computed. Skip-gram treats each context word equally : the models computes the probability for each word of appearing in the context independently of its distance to the center word

4.4 Negative Sampling

Lets take a second to look at the objective function. Note that the summation over $|V|$ is computationally huge! Any update we do or evaluation of the objective function would take $O(|V|)$ time which if we recall is in the millions. A simple idea is we could instead just approximate it.

For every training step, instead of looping over the entire vocabulary, we can just sample several negative examples! We "sample" from a noise distribution ($P_n(w)$) whose probabilities match the ordering of the frequency of the vocabulary. To augment our formulation of the problem to incorporate Negative Sampling, all we need to do is update the:

- objective function

Loss functions J for CBOW and Skip-Gram are expensive to compute because of the softmax normalization, where we sum over all $|V|$ scores!

- gradients
- update rules

MIKOLOV ET AL. present **Negative Sampling** in DISTRIBUTED REPRESENTATIONS OF WORDS AND PHRASES AND THEIR COMPOSITIONALITY. While negative sampling is based on the Skip-Gram model, it is in fact optimizing a different objective. Consider a pair (w, c) of word and context. Did this pair come from the training data? Let's denote by $P(D = 1|w, c)$ the probability that (w, c) came from the corpus data. Correspondingly, $P(D = 0|w, c)$ will be the probability that (w, c) did not come from the corpus data. First, let's model $P(D = 1|w, c)$ with the sigmoid function:

$$P(D = 1|w, c, \theta) = \sigma(v_c^T v_w) = \frac{1}{1 + e^{(-v_c^T v_w)}}$$

Now, we build a new objective function that tries to maximize the probability of a word and context being in the corpus data if it indeed is, and maximize the probability of a word and context not being in the corpus data if it indeed is not. We take a simple maximum likelihood approach of these two probabilities. (Here we take θ to be the parameters of the model, and in our case it is \mathcal{V} and \mathcal{U} .)

$$\begin{aligned} \theta &= \operatorname{argmax}_{\theta} \prod_{(w,c) \in D} P(D = 1|w, c, \theta) \prod_{(w,c) \in \tilde{D}} P(D = 0|w, c, \theta) \\ &= \operatorname{argmax}_{\theta} \prod_{(w,c) \in D} P(D = 1|w, c, \theta) \prod_{(w,c) \in \tilde{D}} (1 - P(D = 1|w, c, \theta)) \\ &= \operatorname{argmax}_{\theta} \sum_{(w,c) \in D} \log P(D = 1|w, c, \theta) + \sum_{(w,c) \in \tilde{D}} \log(1 - P(D = 1|w, c, \theta)) \\ &= \operatorname{argmax}_{\theta} \sum_{(w,c) \in D} \log \frac{1}{1 + \exp(-u_w^T v_c)} + \sum_{(w,c) \in \tilde{D}} \log(1 - \frac{1}{1 + \exp(-u_w^T v_c)}) \\ &= \operatorname{argmax}_{\theta} \sum_{(w,c) \in D} \log \frac{1}{1 + \exp(-u_w^T v_c)} + \sum_{(w,c) \in \tilde{D}} \log(\frac{1}{1 + \exp(u_w^T v_c)}) \end{aligned}$$

Note that maximizing the likelihood is the same as minimizing the negative log likelihood

$$J = - \sum_{(w,c) \in D} \log \frac{1}{1 + \exp(-u_w^T v_c)} - \sum_{(w,c) \in \tilde{D}} \log(\frac{1}{1 + \exp(u_w^T v_c)})$$

Note that \tilde{D} is a "false" or "negative" corpus. Where we would have sentences like "stock boil fish is toy". Unnatural sentences that should get a low probability of ever occurring. We can generate \tilde{D} on the fly by randomly sampling this negative from the word bank.

The **sigmoid** function
 $\sigma(x) = \frac{1}{1+e^{-x}}$
is the 1D version of the softmax and
can be used to model a probability

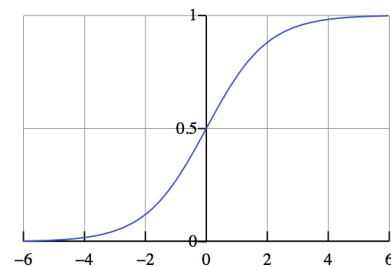


Figure 3: Sigmoid function

For skip-gram, our new objective function for observing the context word $c - m + j$ given the center word c would be

$$-\log \sigma(u_{c-m+j}^T \cdot v_c) - \sum_{k=1}^K \log \sigma(-\tilde{u}_k^T \cdot v_c)$$

For CBOW, our new objective function for observing the center word u_c given the context vector $\hat{v} = \frac{v_{c-m} + v_{c-m+1} + \dots + v_{c+m}}{2m}$ would be

$$-\log \sigma(u_c^T \cdot \hat{v}) - \sum_{k=1}^K \log \sigma(-\tilde{u}_k^T \cdot \hat{v})$$

In the above formulation, $\{\tilde{u}_k | k = 1 \dots K\}$ are sampled from $P_n(w)$. Let's discuss what $P_n(w)$ should be. While there is much discussion of what makes the best approximation, what seems to work best is the Unigram Model raised to the power of $3/4$. Why $3/4$? Here's an example that might help gain some intuition:

$$\begin{aligned} \text{is: } 0.9^{3/4} &= 0.92 \\ \text{Constitution: } 0.09^{3/4} &= 0.16 \\ \text{bombastic: } 0.01^{3/4} &= 0.032 \end{aligned}$$

"Bombastic" is now 3x more likely to be sampled while "is" only went up marginally.

4.5 Hierarchical Softmax

MIKOLOV ET AL. also present hierarchical softmax as a much more efficient alternative to the normal softmax. In practice, hierarchical softmax tends to be better for infrequent words, while negative sampling works better for frequent words and lower dimensional vectors.

Hierarchical softmax uses a binary tree to represent all words in the vocabulary. Each leaf of the tree is a word, and there is a unique path from root to leaf. In this model, there is *no output representation for words*. Instead, each node of the graph (except the root and the leaves) is associated to a vector that the model is going to learn.

In this model, the probability of a word w given a vector w_i , $P(w|w_i)$, is equal to the probability of a random walk starting in the root and ending in the leaf node corresponding to w . The main advantage in computing the probability this way is that the cost is only $O(\log(|V|))$, corresponding to the length of the path.

Let's introduce some notation. Let $L(w)$ be the number of nodes in the path from the root to the leaf w . For instance, $L(w_2)$ in Figure 4 is 3. Let's write $n(w, i)$ as the i -th node on this path with associated

To compare with the regular softmax loss for skip-gram

$$-u_{c-m+j}^T v_c + \log \sum_{k=1}^{|V|} \exp(u_k^T v_c)$$

To compare with the regular softmax loss for CBOW

$$-u_c^T \hat{v} + \log \sum_{j=1}^{|V|} \exp(u_j^T \hat{v})$$

Hierarchical Softmax uses a binary tree where leaves are the words. The probability of a word being the output word is defined as the probability of a random walk from the root to that word's leaf. Computational cost becomes $O(\log(|V|))$ instead of $O(|V|)$.



Figure 4: Binary tree for Hierarchical softmax

vector $v_{n(w,i)}$. So $n(w, 1)$ is the root, while $n(w, L(w))$ is the father of w . Now for each inner node n , we arbitrarily choose one of its children and call it $ch(n)$ (e.g. always the left node). Then, we can compute the probability as

$$P(w|w_i) = \prod_{j=1}^{L(w)-1} \sigma([n(w, j+1) = ch(n(w, j))]) \cdot v_{n(w,j)}^T v_{w_i}$$

where

$$[x] = \begin{cases} 1 & \text{if } x \text{ is true} \\ -1 & \text{otherwise} \end{cases}$$

and $\sigma(\cdot)$ is the sigmoid function.

This formula is fairly dense, so let's examine it more closely.

First, we are computing a product of terms based on the shape of the path from the root ($n(w, 1)$) to the leaf (w). If we assume $ch(n)$ is always the left node of n , then term $[n(w, j+1) = ch(n(w, j))]$ returns 1 when the path goes left, and -1 if right.

Furthermore, the term $[n(w, j+1) = ch(n(w, j))]$ provides normalization. At a node n , if we sum the probabilities for going to the left and right node, you can check that for any value of $v_n^T v_{w_i}$,

$$\sigma(v_n^T v_{w_i}) + \sigma(-v_n^T v_{w_i}) = 1$$

The normalization also ensures that $\sum_{w=1}^{|V|} P(w|w_i) = 1$, just as in the original softmax.

Finally, we compare the similarity of our input vector v_{w_i} to each inner node vector $v_{n(w,j)}$ using a dot product. Let's run through an example. Taking w_2 in Figure 4, we must take two left edges and then a right edge to reach w_2 from the root, so

$$\begin{aligned} P(w_2|w_i) &= p(n(w_2, 1), \text{left}) \cdot p(n(w_2, 2), \text{left}) \cdot p(n(w_2, 3), \text{right}) \\ &= \sigma(v_{n(w_2,1)}^T v_{w_i}) \cdot \sigma(v_{n(w_2,2)}^T v_{w_i}) \cdot \sigma(-v_{n(w_2,3)}^T v_{w_i}) \end{aligned}$$

To train the model, our goal is still to minimize the negative log likelihood $-\log P(w|w_i)$. But instead of updating output vectors per word, we update the vectors of the nodes in the binary tree that are in the path from root to leaf node.

The speed of this method is determined by the way in which the binary tree is constructed and words are assigned to leaf nodes. MIKOLOV ET AL. use a binary Huffman tree, which assigns frequent words shorter paths in the tree.

References

[Bengio et al., 2003] Bengio, Y., Ducharme, R., Vincent, P., and Janvin, C. (2003). A neural probabilistic language model. *J. Mach. Learn. Res.*, 3:1137–1155.

- [Collobert et al., 2011] Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., and Kuksa, P. P. (2011). Natural language processing (almost) from scratch. *CoRR*, abs/1103.0398.
- [Mikolov et al., 2013] Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781.
- [Rong, 2014] Rong, X. (2014). word2vec parameter learning explained. *CoRR*, abs/1411.2738.
- [Rumelhart et al., 1988] Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1988). Neurocomputing: Foundations of research. chapter Learning Representations by Back-propagating Errors, pages 696–699. MIT Press, Cambridge, MA, USA.

CS224n: Natural Language Processing with Deep Learning¹

Lecture Notes: Part II

Word Vectors II: GloVe, Evaluation and Training²

Winter 2019

¹ Course Instructors: Christopher Manning, Richard Socher

² Authors: Rohit Mundra, Emma Peng, Richard Socher, Ajay Sohmshtetty, Amita Kamath

Keyphrases: Global Vectors for Word Representation (GloVe). Intrinsic and extrinsic evaluations. Effect of hyperparameters on analogy evaluation tasks. Correlation of human judgment with word vector distances. Dealing with ambiguity in word using contexts. Window classification.

This set of notes first introduces the GloVe model for training word vectors. Then it extends our discussion of word vectors (interchangeably called word embeddings) by seeing how they can be evaluated intrinsically and extrinsically. As we proceed, we discuss the example of word analogies as an intrinsic evaluation technique and how it can be used to tune word embedding techniques. We then discuss training model weights/parameters and word vectors for extrinsic tasks. Lastly we motivate artificial neural networks as a class of models for natural language processing tasks.

1 Global Vectors for Word Representation (GloVe)³

1.1 Comparison with Previous Methods

So far, we have looked at two main classes of methods to find word embeddings. The first set are count-based and rely on matrix factorization (e.g. LSA, HAL). While these methods effectively leverage global statistical information, they are primarily used to capture word similarities and do poorly on tasks such as word analogy, indicating a sub-optimal vector space structure. The other set of methods are shallow window-based (e.g. the skip-gram and the CBOW models), which learn word embeddings by making predictions in local context windows. These models demonstrate the capacity to capture complex linguistic patterns beyond word similarity, but fail to make use of the global co-occurrence statistics.

In comparison, GloVe consists of a weighted least squares model that trains on global word-word co-occurrence counts and thus makes efficient use of statistics. The model produces a word vector space with meaningful sub-structure. It shows state-of-the-art performance on the word analogy task, and outperforms other current methods on several word similarity tasks.

³ This section is based on the GloVe paper by Pennington et al.:

Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global Vectors for Word Representation

GloVe:

- Using global statistics to predict the probability of word j appearing in the context of word i with a least squares objective

1.2 Co-occurrence Matrix

Let X denote the word-word co-occurrence matrix, where X_{ij} indicates the number of times word j occur in the context of word i . Let $X_i = \sum_k X_{ik}$ be the number of times any word k appears in the context of word i . Finally, let $P_{ij} = P(w_j|w_i) = \frac{X_{ij}}{X_i}$ be the probability of j appearing in the context of word i .

Populating this matrix requires a single pass through the entire corpus to collect the statistics. For large corpora, this pass can be computationally expensive, but it is a one-time up-front cost.

Co-occurrence Matrix:

- X : word-word co-occurrence matrix
- X_{ij} : number of times word j occur in the context of word i
- $X_i = \sum_k X_{ik}$: the number of times any word k appears in the context of word i
- $P_{ij} = P(w_j|w_i) = \frac{X_{ij}}{X_i}$: the probability of j appearing in the context of word i

1.3 Least Squares Objective

Recall that for the skip-gram model, we use softmax to compute the probability of word j appears in the context of word i :

$$Q_{ij} = \frac{\exp(\vec{u}_j^T \vec{v}_i)}{\sum_{w=1}^W \exp(\vec{u}_w^T \vec{v}_i)}$$

Training proceeds in an on-line, stochastic fashion, but the implied global cross-entropy loss can be calculated as:

$$J = - \sum_{i \in \text{corpus}} \sum_{j \in \text{context}(i)} \log Q_{ij}$$

As the same words i and j can appear multiple times in the corpus, it is more efficient to first group together the same values for i and j :

$$J = - \sum_{i=1}^W \sum_{j=1}^W X_{ij} \log Q_{ij}$$

where the value of co-occurring frequency is given by the co-occurrence matrix X . One significant drawback of the cross-entropy loss is that it requires the distribution Q to be properly normalized, which involves the expensive summation over the entire vocabulary. Instead, we use a least square objective in which the normalization factors in P and Q are discarded:

$$\hat{J} = \sum_{i=1}^W \sum_{j=1}^W X_{ij} (\hat{P}_{ij} - \hat{Q}_{ij})^2$$

where $\hat{P}_{ij} = X_{ij}$ and $\hat{Q}_{ij} = \exp(\vec{u}_j^T \vec{v}_i)$ are the unnormalized distributions. This formulation introduces a new problem – X_{ij} often takes on very large values and makes the optimization difficult. An effective change is to minimize the squared error of the logarithms of \hat{P} and \hat{Q} :

$$\begin{aligned}\hat{J} &= \sum_{i=1}^W \sum_{j=1}^W X_i (\log(\hat{P})_{ij} - \log(\hat{Q}_{ij}))^2 \\ &= \sum_{i=1}^W \sum_{j=1}^W X_i (\vec{u}_j^T \vec{v}_i - \log X_{ij})^2\end{aligned}$$

Another observation is that the weighting factor X_i is not guaranteed to be optimal. Instead, we introduce a more general weighting function, which we are free to take to depend on the context word as well:

$$\hat{J} = \sum_{i=1}^W \sum_{j=1}^W f(X_{ij}) (\vec{u}_j^T \vec{v}_i - \log X_{ij})^2$$

1.4 Conclusion

In conclusion, the GloVe model efficiently leverages global statistical information by training only on the nonzero elements in a word-word co-occurrence matrix, and produces a vector space with meaningful sub-structure. It consistently outperforms *word2vec* on the word analogy task, given the same corpus, vocabulary, window size, and training time. It achieves better results faster, and also obtains the best results irrespective of speed.

2 Evaluation of Word Vectors

So far, we have discussed methods such as the *Word2Vec* and *GloVe* methods to train and discover latent vector representations of natural language words in a semantic space. In this section, we discuss how we can quantitatively evaluate the quality of word vectors produced by such techniques.

2.1 Intrinsic Evaluation

Intrinsic evaluation of word vectors is the evaluation of a set of word vectors generated by an embedding technique (such as Word2Vec or GloVe) on specific intermediate subtasks (such as analogy completion). These subtasks are typically simple and fast to compute and thereby allow us to help understand the system used to generate the word vectors. An intrinsic evaluation should typically return to us a number that indicates the performance of those word vectors on the evaluation subtask.

Motivation: Let us consider an example where our final goal is to create a question answering system which uses word vectors

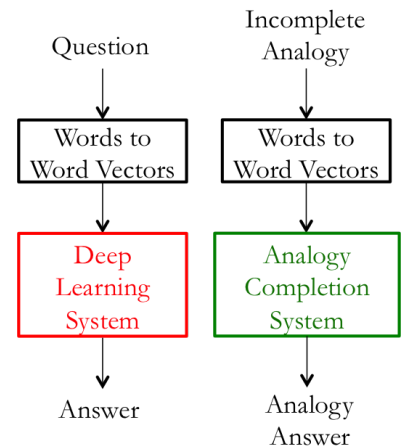


Figure 1: The left subsystem (red) being expensive to train is modified by substituting with a simpler subsystem (green) for intrinsic evaluation.

as inputs. One approach of doing so would be to train a machine learning system that:

1. Takes words as inputs
2. Converts them to word vectors
3. Uses word vectors as inputs for an elaborate machine learning system
4. Maps the output word vectors by this system back to natural language words
5. Produces words as answers

Of course, in the process of making such a state-of-the-art question-answering system, we will need to create optimal word-vector representations since they are used in downstream subsystems (such as deep neural networks). To do this in practice, we will need to tune many hyperparameters in the Word2Vec subsystem (such as the dimension of the word vector representation). While the idealistic approach is to retrain the entire system after any parametric changes in the Word2Vec subsystem, this is impractical from an engineering standpoint because the machine learning system (in step 3) is typically a deep neural network with millions of parameters that takes very long to train. In such a situation, we would want to come up with a simple intrinsic evaluation technique which can provide a measure of "goodness" of the word to word vector subsystem. Obviously, a requirement is that the intrinsic evaluation has a positive correlation with the final task performance.

Intrinsic evaluation:

- Evaluation on a specific, intermediate task
- Fast to compute performance
- Helps understand subsystem
- Needs positive correlation with real task to determine usefulness

2.2 Extrinsic Evaluation

Extrinsic evaluation of word vectors is the evaluation of a set of word vectors generated by an embedding technique on the real task at hand. These tasks are typically elaborate and slow to compute. Using our example from above, the system which allows for the evaluation of answers from questions is the extrinsic evaluation system. Typically, optimizing over an underperforming extrinsic evaluation system does not allow us to determine which specific subsystem is at fault and this motivates the need for intrinsic evaluation.

Extrinsic evaluation:

- Is the evaluation on a real task
- Can be slow to compute performance
- Unclear if subsystem is the problem, other subsystems, or internal interactions
- If replacing subsystem improves performance, the change is likely good

2.3 Intrinsic Evaluation Example: Word Vector Analogies

A popular choice for intrinsic evaluation of word vectors is its performance in completing word vector analogies. In a word vector analogy, we are given an incomplete analogy of the form:

a : b :: c : ?

The intrinsic evaluation system then identifies the word vector which maximizes the cosine similarity:

$$d = \operatorname{argmax}_i \frac{(x_b - x_a + x_c)^T x_i}{\|x_b - x_a + x_c\|}$$

This metric has an intuitive interpretation. Ideally, we want $x_b - x_a = x_d - x_c$ (For instance, queen – king = actress – actor). This implies that we want $x_b - x_a + x_c = x_d$. Thus we identify the vector x_d which maximizes the normalized dot-product between the two word vectors (i.e. cosine similarity).

Using intrinsic evaluation techniques such as word-vector analogies should be handled with care (keeping in mind various aspects of the corpus used for pre-training). For instance, consider analogies of the form:

City 1 : State containing City 1 :: City 2 : State containing City 2

Input	Result Produced
Chicago : Illinois :: Houston	Texas
Chicago : Illinois :: Philadelphia	Pennsylvania
Chicago : Illinois :: Phoenix	Arizona
Chicago : Illinois :: Dallas	Texas
Chicago : Illinois :: Jacksonville	Florida
Chicago : Illinois :: Indianapolis	Indiana
Chicago : Illinois :: Austin	Texas
Chicago : Illinois :: Detroit	Michigan
Chicago : Illinois :: Memphis	Tennessee
Chicago : Illinois :: Boston	Massachusetts

In many cases above, there are multiple cities/towns/villages with the same name across the US. Thus, many states would qualify as the right answer. For instance, there are at least 10 places in the US called Phoenix and thus, Arizona need not be the only correct response. Let us now consider analogies of the form:

Capital City 1 : Country 1 :: Capital City 2 : Country 2

In many of the cases above, the resulting city produced by this task has only been the capital in the recent past. For instance, prior to 1997 the capital of Kazakhstan was Almaty. Thus, we can anticipate other issues if our corpus is dated.

The previous two examples demonstrated semantic testing using word vectors. We can also test syntax using word vector analogies. The following intrinsic evaluation tests the word vectors' ability to capture the notion of superlative adjectives:

Similarly, the intrinsic evaluation shown below tests the word vectors' ability to capture the notion of past tense:

Note: We see that word vectors incorporate meaning in cosine distance between words. They also incorporate meaning in different dimensions: for example, "nokia" might be close to "samsung" in some dimension for the reason that they are both electronics companies, but might be close to "finland" in another dimension for a different reason, that Nokia is a Finnish company.

Note: Interesting results can be seen when word vectors are cast down to 2 dimensions (using PCA) and graphed: similar words cluster together. However, it is important to remember that a significant portion of spatial information is lost during dimensionality reduction; thus complex relations between words as described in the Nokia example above may not emerge.

Table 1: Here are **semantic** word vector analogies (intrinsic evaluation) that may suffer from different cities having the same name

Input	Result Produced
Abuja : Nigeria : : Accra	Ghana
Abuja : Nigeria : : Algiers	Algeria
Abuja : Nigeria : : Amman	Jordan
Abuja : Nigeria : : Ankara	Turkey
Abuja : Nigeria : : Antananarivo	Madagascar
Abuja : Nigeria : : Apia	Samoa
Abuja : Nigeria : : Ashgabat	Turkmenistan
Abuja : Nigeria : : Asmara	Eritrea
Abuja : Nigeria : : Astana	Kazakhstan

Table 2: Here are **semantic** word vector analogies (intrinsic evaluation) that may suffer from countries having different capitals at different points in time

Input	Result Produced
bad : worst : : big	biggest
bad : worst : : bright	brightest
bad : worst : : cold	coldest
bad : worst : : cool	coolest
bad : worst : : dark	darkest
bad : worst : : easy	easiest
bad : worst : : fast	fastest
bad : worst : : good	best
bad : worst : : great	greatest

Table 3: Here are **syntactic** word vector analogies (intrinsic evaluation) that test the notion of superlative adjectives

2.4 Intrinsic Evaluation Tuning Example: Analogy Evaluations

We now explore some of the hyperparameters in word vector embedding techniques (such as Word2Vec and GloVe) that can be tuned using an intrinsic evaluation system (such as an analogy completion system). Let us first see how different methods for creating word-vector embeddings have performed (in recent research work) under the same hyperparameters on an analogy evaluation task:

Some parameters we might consider tuning for a word embedding technique on intrinsic evaluation tasks are:

- Dimension of word vectors
- Corpus size
- Corpus source/type
- Context window size
- Context symmetry

Can you think of other hyperparameters tunable at this stage?

Input	Result Produced
dancing : danced : : decreasing	decreased
dancing : danced : : describing	described
dancing : danced : : enhancing	enhanced
dancing : danced : : falling	fell
dancing : danced : : feeding	fed
dancing : danced : : flying	flew
dancing : danced : : generating	generated
dancing : danced : : going	went
dancing : danced : : hiding	hid
dancing : danced : : hitting	hit

Table 4: Here are **syntactic** word vector analogies (intrinsic evaluation) that test the notion of past tense

Model	Dimension	Size	Semantics	Syntax	Total
ivLBL	100	1.5B	55.9	50.1	53.2
HPCA	100	1.6B	4.2	16.4	10.8
GloVe	100	1.6B	67.5	54.3	60.3
SG	300	1B	61	61	61
CBOW	300	1.6B	16.1	52.6	36.1
vLBL	300	1.5B	54.2	64.8	60.0
ivLBL	300	1.5B	65.2	63.0	64.0
GloVe	300	1.6B	80.8	61.5	70.3
SVD	300	6B	6.3	8.1	7.3
SVD-S	300	6B	36.7	46.6	42.1
SVD-L	300	6B	56.6	63.0	60.1
CBOW	300	6B	63.6	67.4	65.7
SG	300	6B	73.0	66.0	69.1
GloVe	300	6B	77.4	67.0	71.7
CBOW	1000	6B	57.3	68.9	63.7
SG	1000	6B	66.1	65.1	65.6
SVD-L	300	42B	38.4	58.2	49.2
GloVe	300	42B	81.9	69.3	75.0

Table 5: Here we compare the performance of different models under the use of different hyperparameters and datasets

Inspecting the above table, we can make 3 primary observations:

- **Performance is heavily dependent on the model used for word embedding:**

This is an expected result since different methods try embedding words to vectors using fundamentally different properties (such as co-occurrence count, singular vectors, etc.)

- **Performance increases with larger corpus sizes:**

This happens because of the experience an embedding technique gains with more examples it sees. For instance, an analogy completion example will produce incorrect results if it has not encountered the test words previously.

- **Performance is lower for extremely low dimensional word vectors:**

Lower dimensional word vectors are not able to capture the different meanings of the different words in the corpus. This can be viewed as a high bias problem where our model complexity is too low. For instance, let us consider the words "king", "queen", "man", "woman". Intuitively, we would need to use two dimensions such as "gender" and "leadership" to encode these into 2-bit word vectors. Any lower would fail to capture semantic differences between the four words.

Figure 3 demonstrates how accuracy has been shown to improve with larger corpus.

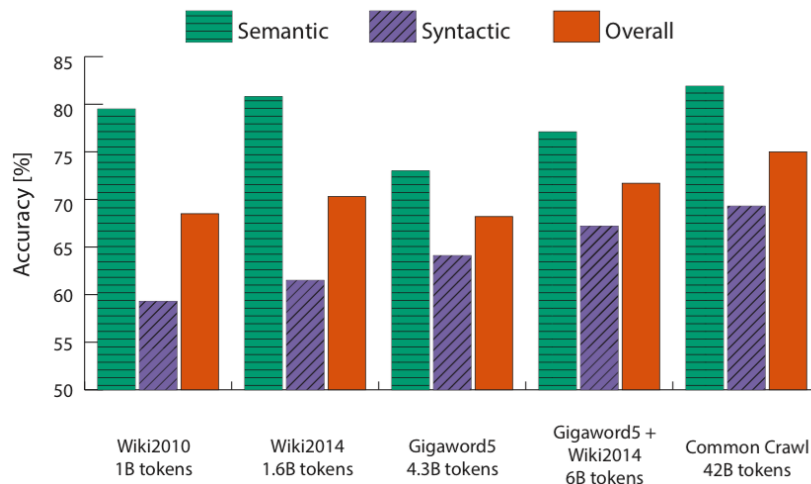


Figure 4 demonstrates how other hyperparameters have been shown to affect the accuracies using GloVe.

Implementation Tip: A window size of 8 around each center word typically works well for GloVe embeddings

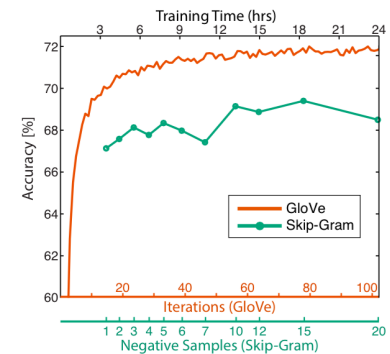


Figure 2: Here we see how training time improves training performance and helps squeeze the last few performance.

Extremely high dimensional vectors: Intuitively, it would seem that these vectors would capture noise in the corpus that doesn't allow generalization, i.e. resulting in high variance, but YIN ET AL. show in ON THE DIMENSIONALITY OF WORD EMBEDDING that skip-gram and GloVe are robust to this over-fitting.

Figure 3: Here we see how performance improves with data size.

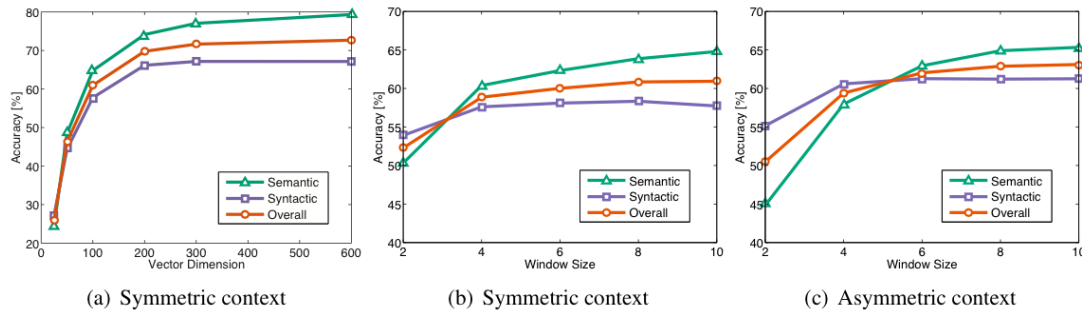


Figure 4: We see how accuracies vary with vector dimension and context window size for GloVe

2.5 Intrinsic Evaluation Example: Correlation Evaluation

Another simple way to evaluate the quality of word vectors is by asking humans to assess the similarity between two words on a fixed scale (say 0-10) and then comparing this with the cosine similarity between the corresponding word vectors. This has been done on various datasets that contain human judgement survey data.

Model	Size	WS353	MC	RG	SCWS	RW
SVD	6B	35.3	35.1	42.5	38.3	25.6
SVD-S	6B	56.5	71.5	71.0	53.6	34.7
SVD-L	6B	65.7	72.7	75.1	56.5	37.0
CBOW	6B	57.2	65.6	68.2	57.0	32.5
SG	6B	62.8	65.2	69.7	58.1	37.2
GloVe	6B	65.8	72.7	77.8	53.9	38.1
SVD-L	42B	74.0	76.4	74.1	58.3	39.9
GloVe	42B	75.9	83.6	82.9	59.6	47.8
CBOW	100B	68.4	79.6	75.4	59.4	45.5

Table 6: Here we see the correlations between of word vector similarities using different embedding techniques with different human judgment datasets

2.6 Further Reading: Dealing With Ambiguity

One might wonder how we handle the situation where we want to capture the same word with different vectors for its different uses in natural language. For instance, "run" is both a noun and a verb and is used and interpreted differently based on the context. IMPROVING WORD REPRESENTATIONS VIA GLOBAL CONTEXT AND MULTIPLE WORD PROTOTYPES (HUANG ET AL, 2012) describes how such cases can also be handled in NLP. The essence of the method is the following:

1. Gather fixed size context windows of all occurrences of the word (for instance, 5 before and 5 after)
2. Each context is represented by a weighted average of the context words' vectors (using idf-weighting)

3. Apply spherical k-means to cluster these context representations.
4. Finally, each word occurrence is re-labeled to its associated cluster and is used to train the word representation for that cluster.

For a more rigorous treatment on this topic, one should refer to the original paper.

3 Training for Extrinsic Tasks

We have so far focused on intrinsic tasks and emphasized their importance in developing a good word embedding technique. Of course, the end goal of most real-world problems is to use the resulting word vectors for some other extrinsic task. Here we discuss the general approach for handling extrinsic tasks.

3.1 Problem Formulation

Most NLP extrinsic tasks can be formulated as classification tasks. For instance, given a sentence, we can classify the sentence to have positive, negative or neutral sentiment. Similarly, in named-entity recognition (NER), given a context and a central word, we want to classify the central word to be one of many classes. For the input, "Jim bought 300 shares of Acme Corp. in 2006", we would like a classified output "[Jim]_{Person} bought 300 shares of [Acme Corp.]_{Organization} in [2006]_{Time}."

For such problems, we typically begin with a training set of the form:

$$\{x^{(i)}, y^{(i)}\}_1^N$$

where $x^{(i)}$ is a d -dimensional word vector generated by some word embedding technique and $y^{(i)}$ is a C -dimensional one-hot vector which indicates the labels we wish to eventually predict (sentiments, other words, named entities, buy/sell decisions, etc.).

In typical machine learning tasks, we usually hold input data and target labels fixed and train weights using optimization techniques (such as gradient descent, L-BFGS, Newton's method, etc.). In NLP applications however, we introduce the idea of retraining the input word vectors when we train for extrinsic tasks. Let us discuss when and why we should consider doing this.

3.2 Retraining Word Vectors

As we have discussed so far, the word vectors we use for extrinsic tasks are initialized by optimizing them over a simpler intrinsic task. In many cases, these pretrained word vectors are a good proxy for

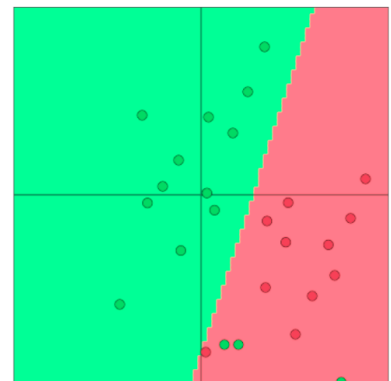


Figure 5: We can classify word vectors using simple linear decision boundaries such as the one shown here (2-D word vectors) using techniques such as logistic regression and SVMs

Implementation Tip: Word vector retraining should be considered for large training datasets. For small datasets, retraining word vectors will likely worsen performance.

optimal word vectors for the extrinsic task and they perform well at the extrinsic task. However, it is also possible that the pretrained word vectors could be trained further (i.e. retrained) using the extrinsic task this time to perform better. However, retraining word vectors can be risky.

If we retrain word vectors using the extrinsic task, we need to ensure that the training set is large enough to cover most words from the vocabulary. This is because Word2Vec or GloVe produce semantically related words to be located in the same part of the word space. When we retrain these words over a small set of the vocabulary, these words are shifted in the word space and as a result, the performance over the final task could actually reduce. Let us explore this idea further using an example. Consider the pretrained vectors to be in a two dimensional space as shown in Figure 6. Here, we see that the word vectors are classified correctly on some extrinsic classification task. Now, if we retrain only two of those vectors because of a limited training set size, then we see in Figure 7 that one of the words gets misclassified because the boundary shifts as a result of word vector updates.

Thus, word vectors should not be retrained if the training data set is small. If the training set is large, retraining may improve performance.

3.3 Softmax Classification and Regularization

Let us consider using the Softmax classification function which has the form:

$$p(y_j = 1|x) = \frac{\exp(W_j \cdot x)}{\sum_{c=1}^C \exp(W_c \cdot x)}$$

Here, we calculate the probability of word vector x being in class j . Using the Cross-entropy loss function, we calculate the loss of such a training example as:

$$-\sum_{j=1}^C y_j \log(p(y_j = 1|x)) = -\sum_{j=1}^C y_j \log\left(\frac{\exp(W_j \cdot x)}{\sum_{c=1}^C \exp(W_c \cdot x)}\right)$$

Of course, the above summation will be a sum over $(C - 1)$ zero values since y_j is 1 only at a single index (at least for now) implying that x belongs to only 1 correct class. Thus, let us define k to be the index of the correct class. Thus, we can now simplify our loss to be:

$$-\log\left(\frac{\exp(W_k \cdot x)}{\sum_{c=1}^C \exp(W_c \cdot x)}\right)$$

We can then extend the above loss to a dataset of N points:

$$-\sum_{i=1}^N \log\left(\frac{\exp(W_k(i) \cdot x^{(i)})}{\sum_{c=1}^C \exp(W_c \cdot x^{(i)})}\right)$$

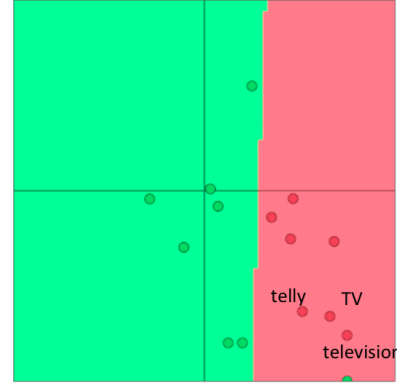


Figure 6: Here, we see that the words "Telly", "TV", and "Television" are classified correctly before retraining. "Telly" and "TV" are present in the extrinsic task training set while "Television" is only present in the test set.

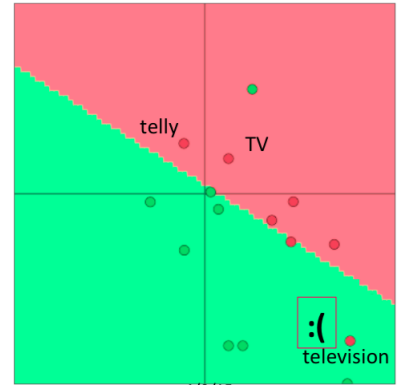


Figure 7: Here, we see that the words "Telly" and "TV" are classified correctly after training, but "Television" is not since it was not present in the training set.

The only difference above is that $k(i)$ is now a function that returns the correct class index for example $x^{(i)}$.

Let us now try to estimate the number of parameters that would be updated if we consider training both, model weights (W), as well word vectors (x). We know that a simple linear decision boundary would require a model that takes in at least one d -dimensional input word vector and produces a distribution over C classes. Thus, to update the model weights, we would be updating $C \cdot d$ parameters. If we update the word vectors for every word in the vocabulary V as well, then we would be updating as many as $|V|$ word vectors, each of which is d -dimensional. Thus, the total number of parameters would be as many as $C \cdot d + |V| \cdot d$ for a simple linear classifier:

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \nabla_{W_{.1}} \\ \vdots \\ \nabla_{W_{.d}} \\ \nabla_{x_{aardvark}} \\ \vdots \\ \nabla_{x_{zebra}} \end{bmatrix}$$

This is an extremely large number of parameters considering how simple the model's decision boundary is - such a large number of parameters is highly prone to overfitting.

To reduce overfitting risk, we introduce a regularization term which poses the Bayesian belief that the parameters (θ) should be small in magnitude (i.e. close to zero):

$$-\sum_{i=1}^N \log \left(\frac{\exp(W_{k(i)} \cdot x^{(i)})}{\sum_{c=1}^C \exp(W_c \cdot x^{(i)})} \right) + \lambda \sum_{k=1}^{C \cdot d + |V| \cdot d} \theta_k^2$$

Minimizing the above cost function reduces the likelihood of the parameters taking on extremely large values just to fit the training set well and may improve generalization if the relative objective weight λ is tuned well. The idea of regularization becomes even more of a requirement once we explore more complex models (such as Neural Networks) which have far more parameters.

3.4 Window Classification

So far we have primarily explored the idea of predicting in extrinsic tasks using a single word vector x . In reality, this is hardly done because of the nature of natural languages. Natural languages tend to use the same word for very different meanings and we typically need to know the context of the word usage to discriminate between meanings. For instance, if you were asked to explain to someone what

... museums in Paris are amazing ...
●●●● ●●●● ●●●● ●●●● ●●●●

Figure 8: Here, we see a central word with a symmetric window of length 2. Such context may help disambiguate between the place Paris and the name Paris.

"to sanction" meant, you would immediately realize that depending on the context "to sanction" could mean "to permit" or "to punish". In most situations, we tend to use a sequence of words as input to the model. A sequence is a central word vector preceded and succeeded by context word vectors. The number of words in the context is also known as the context window size and varies depending on the problem being solved. Generally, narrower window sizes lead to better performance in syntactic tests while wider windows lead to better performance in semantic tests.

In order to modify the previously discussed Softmax model to use windows of words for classification, we would simply substitute $x^{(i)}$ with $x_{window}^{(i)}$ in the following manner:

$$x_{window}^{(i)} = \begin{bmatrix} x^{(i-2)} \\ x^{(i-1)} \\ x^{(i)} \\ x^{(i+1)} \\ x^{(i+2)} \end{bmatrix}$$

As a result, when we evaluate the gradient of the loss with respect to the words, we will receive gradients for the word vectors:

$$\delta_{window} = \begin{bmatrix} \nabla_{x^{(i-2)}} \\ \nabla_{x^{(i-1)}} \\ \nabla_{x^{(i)}} \\ \nabla_{x^{(i+1)}} \\ \nabla_{x^{(i+2)}} \end{bmatrix}$$

The gradient will of course need to be distributed to update the corresponding word vectors in implementation.

3.5 Non-linear Classifiers

We now introduce the need for non-linear classification models such as neural networks. We see in Figure 9 that a linear classifier misclassifies many datapoints. Using a non-linear decision boundary as shown in Figure 10, we manage to classify all training points accurately. Although oversimplified, this is a classic case demonstrating the need for non-linear decision boundaries. In the next set of notes, we study neural networks as a class of non-linear models that have performed particularly well in deep learning applications.

Generally, narrower window sizes lead to better performance in syntactic tests while wider windows lead to better performance in semantic tests.

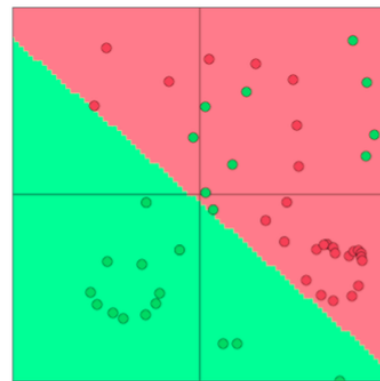


Figure 9: Here, we see that many examples are wrongly classified even though the best linear decision boundary is chosen. This is due linear decision boundaries have limited model capacity for this dataset.

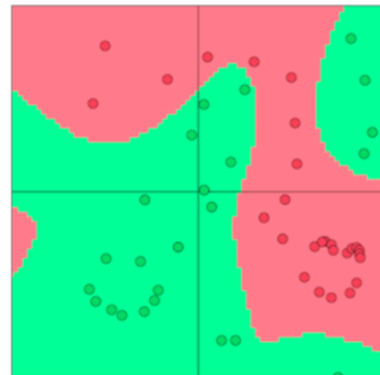


Figure 10: Here, we see that the non-linear decision boundary allows for much better classification of datapoints.

CS224n: Natural Language Processing with Deep Learning ¹

Lecture Notes: Part III

Neural Networks, Backpropagation ²

Winter 2019

¹ Course Instructors: Christopher Manning, Richard Socher

² Authors: Rohit Mundra, Amani Peddada, Richard Socher, Qiaojing Yan

Keyphrases: Neural networks. Forward computation. Backward propagation. Neuron Units. Max-margin Loss. Gradient checks. Xavier parameter initialization. Learning rates. Adagrad.

This set of notes introduces single and multilayer neural networks, and how they can be used for classification purposes. We then discuss how they can be trained using a distributed gradient descent technique known as backpropagation. We will see how the chain rule can be used to make parameter updates sequentially. After a rigorous mathematical discussion of neural networks, we will discuss some practical tips and tricks in training neural networks involving: neuron units (non-linearities), gradient checks, Xavier parameter initialization, learning rates, Adagrad, etc. Lastly, we will motivate the use of recurrent neural networks as a language model.

1 Neural Networks: Foundations

We established in our previous discussions the need for non-linear classifiers since most data are not linearly separable and thus, our classification performance on them is limited. Neural networks are a family of classifiers with non-linear decision boundary as seen in Figure 1. Now that we know the sort of decision boundaries neural networks create, let us see how they manage doing so.

1.1 A Neuron

A neuron is a generic computational unit that takes n inputs and produces a single output. What differentiates the outputs of different neurons is their parameters (also referred to as their weights). One of the most popular choices for neurons is the "sigmoid" or "binary logistic regression" unit. This unit takes an n -dimensional input vector x and produces the scalar activation (output) a . This neuron is also associated with an n -dimensional weight vector, w , and a bias scalar, b . The output of this neuron is then:

$$a = \frac{1}{1 + \exp(-(w^T x + b))}$$

We can also combine the weights and bias term above to equiva-

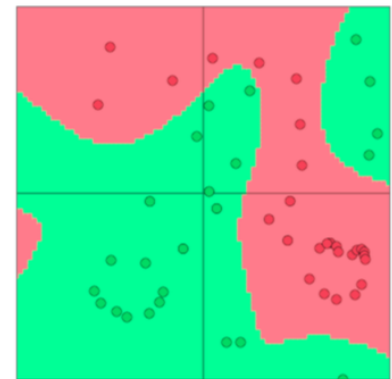


Figure 1: We see here how a non-linear decision boundary separates the data very well. This is the prowess of neural networks.

Fun Fact:

Neural networks are biologically inspired classifiers which is why they are often called "artificial neural networks" to distinguish them from the organic kind. However, in reality human neural networks are so much more capable and complex from artificial neural networks that it is usually better to not draw too many parallels between the two.

Neuron:

A neuron is the fundamental building block of neural networks. We will see that a neuron can be one of many functions that allows for non-linearities to accrue in the network.

lently formulate:

$$a = \frac{1}{1 + \exp(-[w^T \ b] \cdot [x \ 1])}$$

This formulation can be visualized in the manner shown in Figure 2.

1.2 A Single Layer of Neurons

We extend the idea above to multiple neurons by considering the case where the input x is fed as an input to multiple such neurons as shown in Figure 3.

If we refer to the different neurons' weights as $\{w^{(1)}, \dots, w^{(m)}\}$ and the biases as $\{b_1, \dots, b_m\}$, we can say the respective activations are $\{a_1, \dots, a_m\}$:

$$\begin{aligned} a_1 &= \frac{1}{1 + \exp(w^{(1)T}x + b_1)} \\ &\vdots \\ a_m &= \frac{1}{1 + \exp(w^{(m)T}x + b_m)} \end{aligned}$$

Let us define the following abstractions to keep the notation simple and useful for more complex networks:

$$\begin{aligned} \sigma(z) &= \begin{bmatrix} \frac{1}{1+\exp(z_1)} \\ \vdots \\ \frac{1}{1+\exp(z_m)} \end{bmatrix} \\ b &= \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix} \in \mathbb{R}^m \\ W &= \begin{bmatrix} - & w^{(1)T} & - \\ & \dots & \\ - & w^{(m)T} & - \end{bmatrix} \in \mathbb{R}^{m \times n} \end{aligned}$$

We can now write the output of scaling and biases as:

$$z = Wx + b$$

The activations of the sigmoid function can then be written as:

$$\begin{bmatrix} a^{(1)} \\ \vdots \\ a^{(m)} \end{bmatrix} = \sigma(z) = \sigma(Wx + b)$$

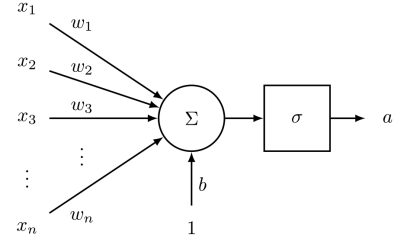


Figure 2: This image captures how in a sigmoid neuron, the input vector x is first scaled, summed, added to a bias unit, and then passed to the squashing sigmoid function.

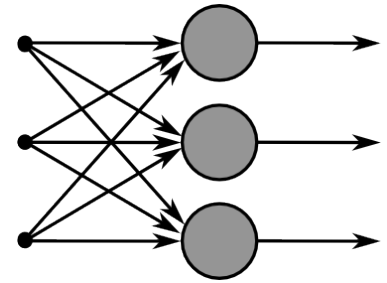


Figure 3: This image captures how multiple sigmoid units are stacked on the right, all of which receive the same input x .

So what do these activations really tell us? Well, one can think of these activations as indicators of the presence of some weighted combination of features. We can then use a combination of these activations to perform classification tasks.

1.3 Feed-forward Computation

So far we have seen how an input vector $x \in \mathbb{R}^n$ can be fed to a layer of sigmoid units to create activations $a \in \mathbb{R}^m$. But what is the intuition behind doing so? Let us consider the following named-entity recognition (NER) problem in NLP as an example:

"Museums in Paris are amazing"

Here, we want to classify whether or not the center word *"Paris"* is a named-entity. In such cases, it is very likely that we would not just want to capture the presence of words in the window of word vectors but some other interactions between the words in order to make the classification. For instance, maybe it should matter that *"Museums"* is the first word only if *"in"* is the second word. Such non-linear decisions can often not be captured by inputs fed directly to a Softmax function but instead require the scoring of the intermediate layer discussed in Section 1.2. We can thus use another matrix $U \in \mathbb{R}^{m \times 1}$ to generate an unnormalized score for a classification task from the activations:

$$s = U^T a = U^T f(Wx + b)$$

where f is the activation function.

Analysis of Dimensions: If we represent each word using a 4-dimensional word vector and we use a 5-word window as input (as in the above example), then the input $x \in \mathbb{R}^{20}$. If we use 8 sigmoid units in the hidden layer and generate 1 score output from the activations, then $W \in \mathbb{R}^{8 \times 20}$, $b \in \mathbb{R}^8$, $U \in \mathbb{R}^{8 \times 1}$, $s \in \mathbb{R}$.

1.4 Maximum Margin Objective Function

Like most machine learning models, neural networks also need an optimization objective, a measure of error or goodness which we want to minimize or maximize respectively. Here, we will discuss a popular error metric known as the maximum margin objective. The idea behind using this objective is to ensure that the score computed for "true" labeled data points is higher than the score computed for "false" labeled data points.

Using the previous example, if we call the score computed for the "true" labeled window *"Museums in Paris are amazing"* as s and the

Dimensions for a single hidden layer neural network: If we represent each word using a 4-dimensional word vector and we use a 5-word window as input, then the input $x \in \mathbb{R}^{20}$. If we use 8 sigmoid units in the hidden layer and generate 1 score output from the activations, then $W \in \mathbb{R}^{8 \times 20}$, $b \in \mathbb{R}^8$, $U \in \mathbb{R}^{8 \times 1}$, $s \in \mathbb{R}$. The stage-wise feed-forward computation is then:

$$\begin{aligned} z &= Wx + b \\ a &= \sigma(z) \\ s &= U^T a \end{aligned}$$

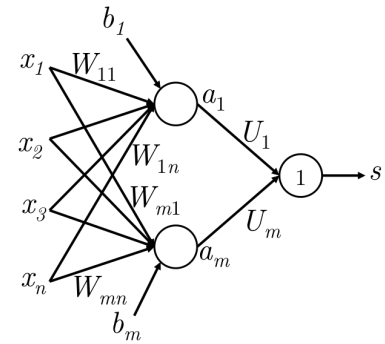


Figure 4: This image captures how a simple feed-forward network might compute its output.

score computed for the "false" labeled window "Not all museums in Paris" as s_c (subscripted as c to signify that the window is "corrupt").

Then, our objective function would be to maximize $(s - s_c)$ or to minimize $(s_c - s)$. However, we modify our objective to ensure that error is only computed if $s_c > s \Rightarrow (s_c - s) > 0$. The intuition behind doing this is that we only care the the "true" data point have a higher score than the "false" data point and that the rest does not matter. Thus, we want our error to be $(s_c - s)$ if $s_c > s$ else 0. Thus, our optimization objective is now:

$$\text{minimize } J = \max(s_c - s, 0)$$

However, the above optimization objective is risky in the sense that it does not attempt to create a margin of safety. We would want the "true" labeled data point to score higher than the "false" labeled data point by some positive margin Δ . In other words, we would want error to be calculated if $(s - s_c < \Delta)$ and not just when $(s - s_c < 0)$. Thus, we modify the optimization objective:

$$\text{minimize } J = \max(\Delta + s_c - s, 0)$$

We can scale this margin such that it is $\Delta = 1$ and let the other parameters in the optimization problem adapt to this without any change in performance. For more information on this, read about functional and geometric margins - a topic often covered in the study of Support Vector Machines. Finally, we define the following optimization objective which we optimize over all training windows:

$$\text{minimize } J = \max(1 + s_c - s, 0)$$

In the above formulation $s_c = U^T f(Wx_c + b)$ and $s = U^T f(Wx + b)$.

1.5 Training with Backpropagation – Elemental

In this section we discuss how we train the different parameters in the model when the cost J discussed in Section 1.4 is positive. No parameter updates are necessary if the cost is 0. Since we typically update parameters using gradient descent (or a variant such as SGD), we typically need the gradient information for any parameter as required in the update equation:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta^{(t)}} J$$

Backpropagation is technique that allows us to use the chain rule of differentiation to calculate loss gradients for any parameter used in the feed-forward computation on the model. To understand this further, let us understand the toy network shown in Figure 5 for which we will perform backpropagation.

The max-margin objective function is most commonly associated with Support Vector Machines (SVMs)

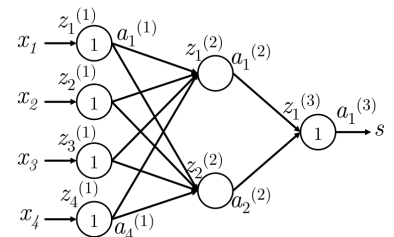


Figure 5: This is a 4-2-1 neural network where neuron j on layer k receives input $z_j^{(k)}$ and produces activation output $a_j^{(k)}$.

Here, we use a neural network with a single hidden layer and a single unit output. Let us establish some **notation** that will make it easier to generalize this model later:

- x_i is an input to the neural network.
- s is the output of the neural network.
- Each layer (including the input and output layers) has neurons which receive an input and produce an output. The j -th neuron of layer k receives the scalar input $z_j^{(k)}$ and produces the scalar activation output $a_j^{(k)}$.
- We will call the backpropagated error calculated at $z_j^{(k)}$ as $\delta_j^{(k)}$.
- Layer 1 refers to the input layer and not the first hidden layer. For the input layer, $x_j = z_j^{(1)} = a_j^{(1)}$.
- $W^{(k)}$ is the transfer matrix that maps the output from the k -th layer to the input to the $(k+1)$ -th. Thus, $W^{(1)} = W$ and $W^{(2)} = U$ to put this new generalized notation in perspective of Section 1.3.

Let us begin: Suppose the cost $J = (1 + s_c - s)$ is positive and we want to perform the update of parameter $W_{14}^{(1)}$ (in Figure 5 and Figure 6), we must realize that $W_{14}^{(1)}$ only contributes to $z_1^{(2)}$ and thus $a_1^{(2)}$. This fact is crucial to understanding backpropagation – backpropagated gradients are only affected by values they contribute to. $a_1^{(2)}$ is consequently used in the forward computation of score by multiplication with $W_1^{(2)}$. We can see from the max-margin loss that:

$$\frac{\partial J}{\partial s} = -\frac{\partial J}{\partial s_c} = -1$$

Therefore we will work with $\frac{\partial s}{\partial W_{ij}^{(1)}}$ here for simplicity. Thus,

$$\begin{aligned} \frac{\partial s}{\partial W_{ij}^{(1)}} &= \frac{\partial W^{(2)} a^{(2)}}{\partial W_{ij}^{(1)}} = \frac{\partial W_i^{(2)} a_i^{(2)}}{\partial W_{ij}^{(1)}} = W_i^{(2)} \frac{\partial a_i^{(2)}}{\partial W_{ij}^{(1)}} \\ \Rightarrow W_i^{(2)} \frac{\partial a_i^{(2)}}{\partial W_{ij}^{(1)}} &= W_i^{(2)} \frac{\partial a_i^{(2)}}{\partial z_i^{(2)}} \frac{\partial z_i^{(2)}}{\partial W_{ij}^{(1)}} \\ &= W_i^{(2)} \frac{f(z_i^{(2)})}{\partial z_i^{(2)}} \frac{\partial z_i^{(2)}}{\partial W_{ij}^{(1)}} \\ &= W_i^{(2)} f'(z_i^{(2)}) \frac{\partial z_i^{(2)}}{\partial W_{ij}^{(1)}} \end{aligned}$$

Backpropagation Notation:

- x_i is an input to the neural network.
- s is the output of the neural network.
- The j -th neuron of layer k receives the scalar input $z_j^{(k)}$ and produces the scalar activation output $a_j^{(k)}$.
- For the input layer, $x_j = z_j^{(1)} = a_j^{(1)}$.
- $W^{(k)}$ is the transfer matrix that maps the output from the k -th layer to the input to the $(k+1)$ -th. Thus, $W^{(1)} = W$ and $W^{(2)} = U^T$ using notation from Section 1.3.

$$\begin{aligned}
 &= W_i^{(2)} f'(z_i^{(2)}) \frac{\partial}{\partial W_{ij}^{(1)}} (b_i^{(1)} + a_1^{(1)} W_{i1}^{(1)} + a_2^{(1)} W_{i2}^{(1)} + a_3^{(1)} W_{i3}^{(1)} + a_4^{(1)} W_{i4}^{(1)}) \\
 &= W_i^{(2)} f'(z_i^{(2)}) \frac{\partial}{\partial W_{ij}^{(1)}} (b_i^{(1)} + \sum_k a_k^{(1)} W_{ik}^{(1)}) \\
 &= W_i^{(2)} f'(z_i^{(2)}) a_j^{(1)} \\
 &= \delta_i^{(2)} \cdot a_j^{(1)}
 \end{aligned}$$

We see above that the gradient reduces to the product $\delta_i^{(2)} \cdot a_j^{(1)}$ where $\delta_i^{(2)}$ is essentially the error propagating backwards from the i -th neuron in layer 2. $a_j^{(1)}$ is an input fed to i -th neuron in layer 2 when scaled by W_{ij} .

Let us discuss the "error sharing/distribution" interpretation of backpropagation better using Figure 6 as an example. Say we were to update $W_{14}^{(1)}$:

1. We start with the an error signal of 1 propagating backwards from $a_1^{(3)}$.
2. We then multiply this error by the local gradient of the neuron which maps $z_1^{(3)}$ to $a_1^{(3)}$. This happens to be 1 in this case and thus, the error is still 1. This is now known as $\delta_1^{(3)} = 1$.
3. At this point, the error signal of 1 has reached $z_1^{(3)}$. We now need to distribute the error signal so that the "fair share" of the error reaches to $a_1^{(2)}$.
4. This amount is the (error signal at $z_1^{(3)} = \delta_1^{(3)} \times W_1^{(2)} = W_1^{(2)}$. Thus, the error at $a_1^{(2)} = W_1^{(2)}$.
5. As we did in step 2, we need to move the error across the neuron which maps $z_1^{(2)}$ to $a_1^{(2)}$. We do this by multiplying the error signal at $a_1^{(2)}$ by the local gradient of the neuron which happens to be $f'(z_1^{(2)})$.
6. Thus, the error signal at $z_1^{(2)}$ is $f'(z_1^{(2)}) W_1^{(2)}$. This is known as $\delta_1^{(2)}$.
7. Finally, we need to distribute the "fair share" of the error to $W_{14}^{(1)}$ by simply multiplying it by the input it was responsible for forwarding, which happens to be $a_4^{(1)}$.
8. Thus, the gradient of the loss with respect to $W_{14}^{(1)}$ is calculated to be $a_4^{(1)} f'(z_1^{(2)}) W_1^{(2)}$.

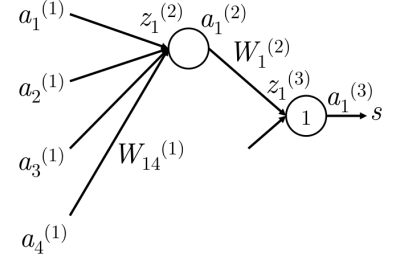


Figure 6: This subnetwork shows the relevant parts of the network required to update $W_{ij}^{(1)}$

Notice that the result we arrive at using this approach is exactly the same as that we arrived at using explicit differentiation earlier. Thus, we can calculate error gradients with respect to a parameter in the network using either the chain rule of differentiation or using an error sharing and distributed flow approach – both of these approaches happen to do the exact same thing but it might be helpful to think about them one way or another.

Bias Updates: Bias terms (such as $b_1^{(1)}$) are mathematically equivalent to other weights contributing to the neuron input ($z_1^{(2)}$) as long as the input being forwarded is 1. As such, the bias gradients for neuron i on layer k is simply $\delta_i^{(k)}$. For instance, if we were updating $b_1^{(1)}$ instead of $W_{14}^{(1)}$ above, the gradient would simply be $f'(z_1^{(2)})W_1^{(2)}$.

Generalized steps to propagate $\delta^{(k)}$ to $\delta^{(k-1)}$:

1. We have error $\delta_i^{(k)}$ propagating backwards from $z_i^{(k)}$, i.e. neuron i at layer k . See Figure 7.
2. We propagate this error backwards to $a_j^{(k-1)}$ by multiplying $\delta_i^{(k)}$ by the path weight $W_{ij}^{(k-1)}$.
3. Thus, the error received at $a_j^{(k-1)}$ is $\delta_i^{(k)} W_{ij}^{(k-1)}$.
4. However, $a_j^{(k-1)}$ may have been forwarded to multiple nodes in the next layer as shown in Figure 8. It should receive responsibility for errors propagating backward from node m in layer k too, using the exact same mechanism.
5. Thus, error received at $a_j^{(k-1)}$ is $\delta_i^{(k)} W_{ij}^{(k-1)} + \delta_m^{(k)} W_{mj}^{(k-1)}$.
6. In fact, we can generalize this to be $\sum_i \delta_i^{(k)} W_{ij}^{(k-1)}$.
7. Now that we have the correct error at $a_j^{(k-1)}$, we move it across neuron j at layer $k-1$ by multiplying with the local gradient $f'(z_j^{(k-1)})$.
8. Thus, the error that reaches $z_j^{(k-1)}$, called $\delta_j^{(k-1)}$ is $f'(z_j^{(k-1)}) \sum_i \delta_i^{(k)} W_{ij}^{(k-1)}$.

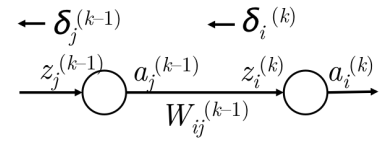


Figure 7: Propagating error from $\delta^{(k)}$ to $\delta^{(k-1)}$

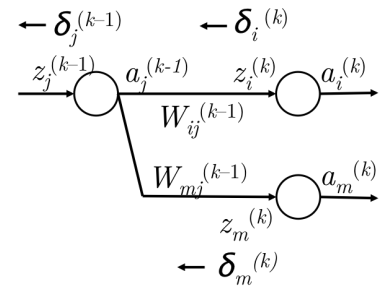


Figure 8: Propagating error from $\delta^{(k)}$ to $\delta^{(k-1)}$

1.6 Training with Backpropagation – Vectorized

So far, we discussed how to calculate gradients for a given parameter in the model. Here we will generalize the approach above so that we update weight matrices and bias vectors all at once. Note that these are simply extensions of the above model that will help build intuition for the way error propagation can be done at a matrix-vector level.

For a given parameter $W_{ij}^{(k)}$, we identified that the error gradient is simply $\delta_i^{(k+1)} \cdot a_j^{(k)}$. As a reminder, $W^{(k)}$ is the matrix that maps $a^{(k)}$ to $z^{(k+1)}$. We can thus establish that the error gradient for the entire matrix $W^{(k)}$ is:

$$\nabla_{W^{(k)}} = \begin{bmatrix} \delta_1^{(k+1)} a_1^{(k)} & \delta_1^{(k+1)} a_2^{(k)} & \dots \\ \delta_2^{(k+1)} a_1^{(k)} & \delta_2^{(k+1)} a_2^{(k)} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} = \delta^{(k+1)} a^{(k)T}$$

Thus, we can write an entire matrix gradient using the outer product of the error vector propagating into the matrix and the activations forwarded by the matrix.

Now, we will see how we can calculate the error vector $\delta^{(k)}$. We established earlier using Figure 8 that $\delta_j^{(k)} = f'(z_j^{(k)}) \sum_i \delta_i^{(k+1)} W_{ij}^{(k)}$. This can easily generalize to matrices such that:

$$\delta^{(k)} = f'(z^{(k)}) \circ (W^{(k)T} \delta^{(k+1)})$$

In the above formulation, the \circ operator corresponds to an element wise product between elements of vectors ($\circ : \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}^N$).

Error propagates from layer $(k+1)$ to (k) in the following manner:

$$\delta^{(k)} = f'(z^{(k)}) \circ (W^{(k)T} \delta^{(k+1)})$$

Of course, this assumes that in the forward propagation the signal $z^{(k)}$ first goes through activation neurons f to generate activations $a^{(k)}$ and are then linearly combined to yield $z^{(k+1)}$ via transfer matrix $W^{(k)}$.

Computational efficiency: Having explored element-wise updates as well as vector-wise updates, we must realize that the vectorized implementations run substantially faster in scientific computing environments such as MATLAB or Python (using NumPy/SciPy packages). Thus, we should use vectorized implementation in practice. Furthermore, we should also reduce redundant calculations in backpropagation - for instance, notice that $\delta^{(k)}$ depends directly on $\delta^{(k+1)}$. Thus, we should ensure that when we update $W^{(k)}$ using $\delta^{(k+1)}$, we save $\delta^{(k+1)}$ to later derive $\delta^{(k)}$ – and we then repeat this for $(k-1) \dots (1)$. Such a recursive procedure is what makes backpropagation a computationally affordable procedure.

2 Neural Networks: Tips and Tricks

Having discussed the mathematical foundations of neural networks, we will now dive into some tips and tricks commonly employed

when using neural networks in practice.

2.1 Gradient Check

In the last section, we discussed in detail how to calculate error gradients/updates for parameters in a neural network model via calculus-based (analytic) methods. Here we now introduce a technique of *numerically* approximating these gradients – though too computationally inefficient to be used directly for training the networks, this method will allow us to very precisely estimate the derivative with respect to any parameter; it can thus serve as a useful sanity check on the correctness of our analytic derivatives. Given a model with parameter vector θ and loss function J , the numerical gradient around θ_i is simply given by **centered difference formula**:

$$f'(\theta) \approx \frac{J(\theta^{(i+)}) - J(\theta^{(i-)})}{2\epsilon}$$

where ϵ is a small number (usually around $1e^{-5}$). The term $J(\theta^{(i+)})$ is simply the error calculated on a forward pass for a given input when we perturb the parameter θ 's i^{th} element by $+\epsilon$. Similarly, the term $J(\theta^{(i-)})$ is the error calculated on a forward pass for the same input when we perturb the parameter θ 's i^{th} element by $-\epsilon$. Thus, using two forward passes, we can approximate the gradient with respect to any given parameter element in the model. We note that this definition of the numerical gradient follows very naturally from the definition of the derivative, where, in the scalar case,

$$f'(x) \approx \frac{f(x + \epsilon) - f(x)}{\epsilon}$$

Of course, there is a slight difference – the definition above only perturbs x in the positive direction to compute the gradient. While it would have been perfectly acceptable to define the numerical gradient in this way, in practice it is often more precise and stable to use the centered difference formula, where we perturb a parameter in both directions. The intuition is that to get a better approximation of the derivative/slope around a point, we need to examine the function f 's behavior both to the left and right of that point. It can also be shown using Taylor's theorem that the centered difference formula has an error proportional to ϵ^2 , which is quite small, whereas the derivative definition is more error-prone.

Now, a natural question you might ask is, if this method is so precise, why do we not use it to compute all of our network gradients instead of applying back-propagation? The simple answer, as hinted earlier, is inefficiency – recall that every time we want to compute the gradient with respect to an element, we need to make two forward

Gradient checks are a great way to compare analytical and numerical gradients. Analytical gradients should be close and numerical gradients can be calculated using:

$$f'(\theta) \approx \frac{J(\theta^{(i+)}) - J(\theta^{(i-)})}{2\epsilon}$$

$J(\theta^{(i+)})$ and $J(\theta^{(i-)})$ can be evaluated using two forward passes. An implementation of this can be seen in Snippet 2.1.

passes through the network, which will be computationally expensive. Furthermore, many large-scale neural networks can contain millions of parameters, and computing two passes per parameter is clearly not optimal. And, since in optimization techniques such as SGD, we must compute the gradients once per iteration for several thousands of iterations, it is obvious that this method quickly grows intractable. This inefficiency is why we only use gradient check to verify the correctness of our analytic gradients, which are much quicker to compute. A standard implementation of gradient check is shown below:

Snippet 2.1

```
def eval_numerical_gradient(f, x):
    """
    a naive implementation of numerical gradient of f at x
    - f should be a function that takes a single argument
    - x is the point (numpy array) to evaluate the gradient
    at
    """

    fx = f(x) # evaluate function value at original point
    grad = np.zeros(x.shape)
    h = 0.00001

    # iterate over all indexes in x
    it = np.nditer(x, flags=['multi_index'],
                    op_flags=['readwrite'])
    while not it.finished:

        # evaluate function at x+h
        ix = it.multi_index
        old_value = x[ix]
        x[ix] = old_value + h # increment by h
        fxh_left = f(x) # evaluate f(x + h)
        x[ix] = old_value - h # decrement by h
        fxh_right = f(x) # evaluate f(x - h)
        x[ix] = old_value # restore to previous value (very
                           important!)

        # compute the partial derivative
        grad[ix] = (fxh_left - fxh_right) / (2*h) # the slope
        it.iternext() # step to next dimension
    return grad
```

2.2 Regularization

As with many machine learning models, neural networks are highly prone to overfitting, where a model is able to obtain near perfect performance on the training dataset, but loses the ability to generalize to unseen data. A common technique used to address overfitting (an issue also known as the “high-variance problem”) is the incorporation of an L_2 regularization penalty. The idea is that we will simply append an extra term to our loss function J , so that the overall cost is now calculated as:

$$J_R = J + \lambda \sum_{i=1}^L \|W^{(i)}\|_F$$

In the above formulation, $\|W^{(i)}\|_F$ is the Frobenius norm of the matrix $W^{(i)}$ (the i -th weight matrix in the network) and λ is the hyper-parameter controlling how much weight the regularization term has relative to the original cost function. Since we are trying to minimize J_R , what regularization is essentially doing is penalizing weights for being too large while optimizing over the original cost function. Due to the quadratic nature of the Frobenius norm (which computes the sum of the squared elements of a matrix), L_2 -regularization effectively reduces the flexibility of the model and thereby reduces the overfitting phenomenon. Imposing such a constraint can also be interpreted as the prior Bayesian belief that the optimal weights are close to zero – how close depends on the value of λ . Choosing the right value of λ is critical, and must be chosen via hyperparameter-tuning. Too high a value of λ causes most of the weights to be set too close to 0, and the model does not learn anything meaningful from the training data, often obtaining poor accuracy on training, validation, and testing sets. Too low a value, and we fall into the domain of overfitting once again. It must be noted that the bias terms are not regularized and do not contribute to the cost term above – try thinking about why this is the case!

There are indeed other types of regularization that are sometimes used, such as L_1 regularization, which sums over the absolute values (rather than squares) of parameter elements – however, this is less commonly applied in practice since it leads to sparsity of parameter weights. In the next section, we discuss *dropout*, which effectively acts as another form of regularization by randomly dropping (i.e. setting to zero) neurons in the forward pass.

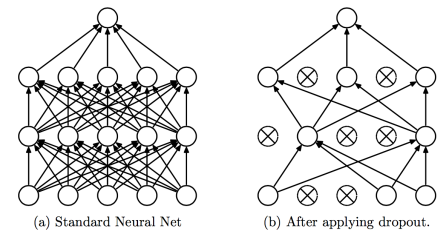
The **Frobenius Norm** of a matrix U is defined as follows:

$$\|U\|_F = \sqrt{\sum_i \sum_j u_{ij}^2}$$

2.3 Dropout

Dropout is a powerful technique for regularization, first introduced by Srivastava et al. in *Dropout: A Simple Way to Prevent Neural Networks from Overfitting*. The idea is simple yet effective – during training, we will randomly “drop” with some probability $(1 - p)$ a subset of neurons during each forward/backward pass (or equivalently, we will keep alive each neuron with a probability p). Then, during testing, we will use the full network to compute our predictions. The result is that the network typically learns more meaningful information from the data, is less likely to overfit, and usually obtains higher performance overall on the task at hand. One intuitive reason why this technique should be so effective is that what dropout is doing is essentially doing is training exponentially many smaller networks at once and averaging over their predictions.

In practice, the way we introduce dropout is that we take the output h of each layer of neurons, and keep each neuron with probability p , and else set it to 0. Then, during back-propagation, we only pass gradients through neurons that were kept alive during the forward pass. Finally, during testing, we compute the forward pass using *all* of the neurons in the network. However, a key subtlety is that in order for dropout to work effectively, the expected output of a neuron during testing should be approximately the same as it was during training – else the magnitude of the outputs could be radically different, and the behavior of the network is no longer well-defined. Thus, we must typically divide the outputs of each neuron during testing by a certain value – it is left as an exercise to the reader to determine what this value should be in order for the expected outputs during training and testing to be equivalent.



Dropout applied to an artificial neural network. Image credits to Srivastava et al.

2.4 Neuron Units

So far we have discussed neural networks that contain sigmoidal neurons to introduce nonlinearities; however in many applications better networks can be designed using other activation functions. Some common choices are listed here with their function and gradient definitions and these can be substituted with the sigmoidal functions discussed above.

Sigmoid: This is the default choice we have discussed; the activation function σ is given by:

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

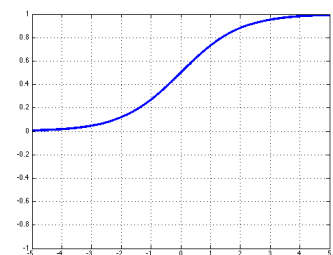


Figure 9: The response of a sigmoid nonlinearity

where $\sigma(z) \in (0, 1)$

The gradient of $\sigma(z)$ is:

$$\sigma'(z) = \frac{-\exp(-z)}{1 + \exp(-z)} = \sigma(z)(1 - \sigma(z))$$

Tanh: The tanh function is an alternative to the sigmoid function that is often found to converge faster in practice. The primary difference between tanh and sigmoid is that tanh output ranges from -1 to 1 while the sigmoid ranges from 0 to 1 .

$$\tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)} = 2\sigma(2z) - 1$$

where $\tanh(z) \in (-1, 1)$

The gradient of $\tanh(z)$ is:

$$\tanh'(z) = 1 - \left(\frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)} \right)^2 = 1 - \tanh^2(z)$$

Hard tanh: The hard tanh function is sometimes preferred over the tanh function since it is computationally cheaper. It does however saturate for magnitudes of z greater than 1 . The activation of the hard tanh is:

$$\text{hardtanh}(z) = \begin{cases} -1 & : z < -1 \\ z & : -1 \leq z \leq 1 \\ 1 & : z > 1 \end{cases}$$

The derivative can also be expressed in a piecewise functional form:

$$\text{hardtanh}'(z) = \begin{cases} 1 & : -1 \leq z \leq 1 \\ 0 & : \text{otherwise} \end{cases}$$

Soft sign: The soft sign function is another nonlinearity which can be considered an alternative to tanh since it too does not saturate as easily as hard clipped functions:

$$\text{softsign}(z) = \frac{z}{1 + |z|}$$

The derivative is the expressed as:

$$\text{softsign}'(z) = \frac{\text{sgn}(z)}{(1 + |z|)^2}$$

where sgn is the signum function which returns ± 1 depending on the sign of z

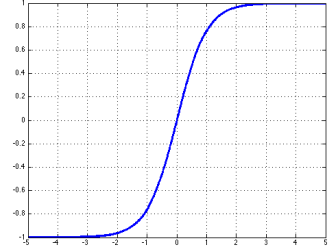


Figure 10: The response of a tanh nonlinearity

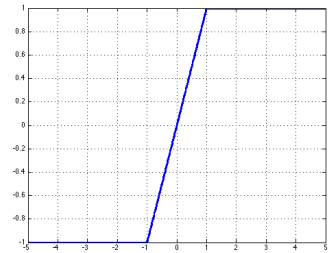


Figure 11: The response of a hard tanh nonlinearity

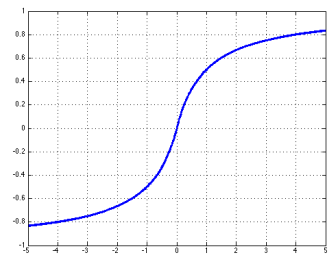


Figure 12: The response of a soft sign nonlinearity

ReLU: The ReLU (Rectified Linear Unit) function is a popular choice of activation since it does not saturate even for larger values of z and has found much success in computer vision applications:

$$\text{rect}(z) = \max(z, 0)$$

The derivative is then the piecewise function:

$$\text{rect}'(z) = \begin{cases} 1 & : z > 0 \\ 0 & : \text{otherwise} \end{cases}$$

Leaky ReLU: Traditional ReLU units by design do not propagate any error for non-positive z – the leaky ReLU modifies this such that a small error is allowed to propagate backwards even when z is negative:

$$\text{leaky}(z) = \max(z, k \cdot z)$$

$$\text{where } 0 < k < 1$$

This way, the derivative is representable as:

$$\text{leaky}'(z) = \begin{cases} 1 & : z > 0 \\ k & : \text{otherwise} \end{cases}$$

2.5 Data Preprocessing

As is the case with machine learning models generally, a key step to ensuring that your model obtains reasonable performance on the task at hand is to perform basic preprocessing on your data. Some common techniques are outlined below.

Mean Subtraction

Given a set of input data X , it is customary to zero-center the data by subtracting the mean feature vector of X from X . An important point is that in practice, the mean is calculated only across the training set, and this mean is subtracted from the training, validation, and testing sets.

Normalization

Another frequently used technique (though perhaps less so than mean subtraction) is to scale every input feature dimension to have similar ranges of magnitudes. This is useful since input features are often measured in different “units”, but we often want to initially consider all features as equally important. The way we accomplish this is by simply dividing the features by their respective standard deviation calculated across the training set.

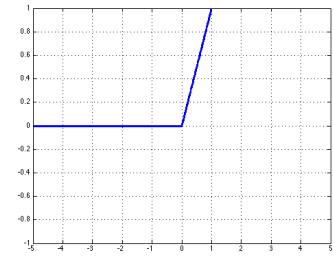


Figure 13: The response of a ReLU nonlinearity

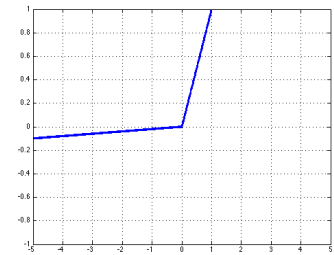


Figure 14: The response of a leaky ReLU nonlinearity

Whitening

Not as commonly used as mean-subtraction + normalization, whitening essentially converts the data to have an identity covariance matrix – that is, features become uncorrelated and have a variance of 1. This is done by first mean-subtracting the data, as usual, to get X' . We can then take the Singular Value Decomposition (SVD) of X' to get matrices U, S, V . We then compute UX' to project X' into the basis defined by the columns of U . We finally divide each dimension of the result by the corresponding singular value in S to scale our data appropriately (if a singular value is zero, we can just divide by a small number instead).

2.6 Parameter Initialization

A key step towards achieving superlative performance with a neural network is initializing the parameters in a reasonable way. A good starting strategy is to initialize the weights to small random numbers normally distributed around 0 – and in practice, this often works acceptably well. However, in UNDERSTANDING THE DIFFICULTY OF TRAINING DEEP FEEDFORWARD NEURAL NETWORKS (2010), XAVIER ET AL study the effect of different weight and bias initialization schemes on training dynamics. The empirical findings suggest that for sigmoid and tanh activation units, faster convergence and lower error rates are achieved when the weights of a matrix $W \in \mathbb{R}^{n^{(l+1)} \times n^{(l)}}$ are initialized randomly with a uniform distribution as follows:

$$W \sim U \left[-\sqrt{\frac{6}{n^{(l)} + n^{(l+1)}}}, \sqrt{\frac{6}{n^{(l)} + n^{(l+1)}}} \right]$$

Where $n^{(l)}$ is the number of input units to W (fan-in) and $n^{(l+1)}$ is the number of output units from W (fan-out). In this parameter initialization scheme, bias units are initialized to 0. This approach attempts to maintain activation variances as well as backpropagated gradient variances across layers. Without such initialization, the gradient variances (which are a proxy for information) generally decrease with backpropagation across layers.

2.7 Learning Strategies

The rate/magnitude of model parameter updates during training can be controlled using the learning rate. In the following naive Gradient Descent formulation, α is the learning rate:

$$\theta^{\text{new}} = \theta^{\text{old}} - \alpha \nabla_{\theta} J_t(\theta)$$

You might think that for fast convergence rates, we should set α to larger values – however faster convergence is not guaranteed with larger convergence rates. In fact, with very large learning rates, we might experience that the loss function actually diverges because the parameters update causes the model to overshoot the convex minima as shown in Figure 15. In non-convex models (most of those we work with), the outcome of a large learning rate is unpredictable, but the chances of diverging loss functions are very high.

The simple solution to avoiding a diverging loss is to use a very small learning rate so that we carefully scan the parameter space – of course, if we use too small a learning rate, we might not converge in a reasonable amount of time, or might get caught in local minima. Thus, as with any other hyperparameter, the learning rate must be tuned effectively.

Since training is the most expensive phase in a deep learning system, some research has attempted to improve this naive approach to setting learning rates. For instance, RONAN COLLOBERT scales the learning rate of a weight W_{ij} (where $W \in \mathbb{R}^{n^{(l+1)} \times n^{(l)}}$) by the inverse square root of the fan-in of the neuron ($n^{(l)}$).

There are several other techniques that have proven to be effective as well – one such method is **annealing**, where, after several iterations, the learning rate is reduced in some way – this method ensures that we start off with a high learning rate and approach a minimum quickly; as we get closer to the minimum, we start lowering our learning rate so that we can find the optimum under a more fine-grained scope. A common way to perform annealing is to reduce the learning rate α by a factor x after every n iterations of learning. Exponential decay is also common, where, the learning rate α at iteration t is given by $\alpha(t) = \alpha_0 e^{-kt}$, where α_0 is the initial learning rate, and k is a hyperparameter. Another approach is to allow the learning rate to decrease over time such that:

$$\alpha(t) = \frac{\alpha_0 \tau}{\max(t, \tau)}$$

In the above scheme, α_0 is a tunable parameter and represents the starting learning rate. τ is also a tunable parameter and represents the time at which the learning rate should start reducing. In practice, this method has been found to work quite well. In the next section we discuss another method for adaptive gradient descent which does not require hand-set learning rates.

2.8 Momentum Updates

Momentum methods, a variant of gradient descent inspired by the study of dynamics and motion in physics, attempt to use the “veloc-

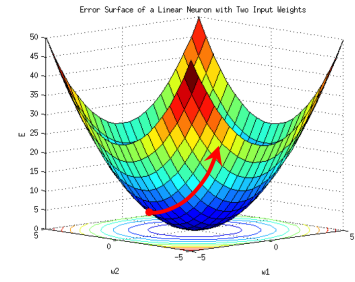


Figure 15: Here we see that updating parameter w_2 with a large learning rate can lead to divergence of the error.

ity” of updates as a more effective update scheme. Pseudocode for momentum updates is shown below:

Snippet 2.2

```
# Computes a standard momentum update
# on parameters x
v = mu*v - alpha*grad_x
x += v
```

2.9 Adaptive Optimization Methods

AdaGrad is an implementation of standard stochastic gradient descent (SGD) with one key difference: the learning rate can vary for each parameter. The learning rate for each parameter depends on the history of gradient updates of that parameter in a way such that parameters with a scarce history of updates are updated faster using a larger learning rate. In other words, parameters that have not been updated much in the past are likelier to have higher learning rates now. Formally:

$$\theta_{t,i} = \theta_{t-1,i} - \frac{\alpha}{\sqrt{\sum_{\tau=1}^t g_{\tau,i}^2}} g_{t,i} \text{ where } g_{t,i} = \frac{\partial}{\partial \theta_i^t} J_t(\theta)$$

In this technique, we see that if the RMS of the history of gradients is extremely low, the learning rate is very high. A simple implementation of this technique is:

Snippet 2.3

```
# Assume the gradient dx and parameter vector x
cache += dx**2
x += - learning_rate * dx / np.sqrt(cache + 1e-8)
```

Other common adaptive methods are RMSProp and Adam, whose update rules are shown below (courtesy of Andrej Karpathy):

Snippet 2.4

```
# Update rule for RMS prop
cache = decay_rate * cache + (1 - decay_rate) * dx**2
x += - learning_rate * dx / (np.sqrt(cache) + eps)
```

Snippet 2.5

```
# Update rule for Adam
```

```
m = beta1*m + (1-beta1)*dx
v = beta2*v + (1-beta2)*(dx**2)
x += - learning_rate * m / (np.sqrt(v) + eps)
```

RMSProp is a variant of AdaGrad that utilizes a moving average of squared gradients – in particular, unlike AdaGrad, its updates do not become monotonically smaller. The Adam update rule is in turn a variant of RMSProp, but with the addition of momentum-like updates. We refer the reader to the respective sources of these methods for more detailed analyses of their behavior.

CS224n: Natural Language Processing with Deep Learning¹

Lecture Notes: Part IV

Dependency Parsing²

Winter 2019

¹ Course Instructors: Christopher Manning, Richard Socher

² Authors: Lisa Wang, Juhi Naik, and Shayne Longpre

Keyphrases: Dependency Parsing.

1 Dependency Grammar and Dependency Structure

Parse trees in NLP, analogous to those in compilers, are used to analyze the syntactic structure of sentences. There are two main types of structures used - constituency structures and dependency structures.

Constituency Grammar uses phrase structure grammar to organize words into nested constituents. This will be covered in more detail in following chapters. We now focus on Dependency Parsing.

Dependency structure of sentences shows which words depend on (modify or are arguments of) which other words. These binary asymmetric relations between the words are called dependencies and are depicted as arrows going from the **head** (or governor, superior, regent) to the **dependent** (or modifier, inferior, subordinate). Usually these dependencies form a tree structure. They are often typed with the name of grammatical relations (subject, prepositional object, apposition, etc.). An example of such a dependency tree is shown in Figure 1. Sometimes a fake ROOT node is added as the head to the whole tree so that every word is a dependent of exactly one node.

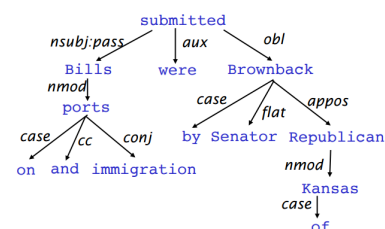


Figure 1: Dependency tree for the sentence "Bills on ports and immigration were submitted by Senator Brownback, Republican of Kansas"

1.1 Dependency Parsing

Dependency parsing is the task of analyzing the syntactic dependency structure of a given input sentence S . The output of a dependency parser is a dependency tree where the words of the input sentence are connected by typed dependency relations. Formally, the dependency parsing problem asks to create a mapping from the input sentence with words $S = w_0w_1...w_n$ (where w_0 is the ROOT) to its dependency tree graph G . Many different variations of dependency-based methods have been developed in recent years, including neural network-based methods, which we will describe later.

To be precise, there are two subproblems in dependency parsing (adapted from Kuebler et al., chapter 1.2):

1. *Learning*: Given a training set D of sentences annotated with dependency graphs, induce a parsing model M that can be used to parse new sentences.

2. *Parsing*: Given a parsing model M and a sentence S , derive the optimal dependency graph D for S according to M .

1.2 Transition-Based Dependency Parsing

Transition-based dependency parsing relies on a state machine which defines the possible transitions to create the mapping from the input sentence to the dependency tree. The *learning problem* is to induce a model which can predict the next transition in the state machine based on the transition history. The *parsing problem* is to construct the optimal sequence of transitions for the input sentence, given the previously induced model. Most transition-based systems do not make use of a formal grammar.

1.3 Greedy Deterministic Transition-Based Parsing

This system was introduced by Nivre in 2003 and was radically different from other methods in use at that time.

This transition system is a state machine, which consists of *states* and *transitions* between those states. The model induces a sequence of transitions from some *initial* state to one of several *terminal* states.

States:

For any sentence $S = w_0w_1...w_n$, a state can be described with a triple $c = (\sigma, \beta, A)$:

1. a stack σ of words w_i from S ,
2. a buffer β of words w_i from S ,
3. a set of dependency arcs A of the form (w_i, r, w_j) , where w_i, w_j are from S , and r describes a dependency relation.

It follows that for any sentence $S = w_0w_1...w_n$,

1. an *initial* state c_0 is of the form $([w_0]_\sigma, [w_1, ..., w_n]_\beta, \emptyset)$ (only the ROOT is on the stack σ , all other words are in the buffer β and no actions have been chosen yet),
2. a *terminal* state has the form $(\sigma, [], A)$.

Transitions:

There are three types of transitions between states:

1. **SHIFT**: Remove the first word in the buffer and push it on top of the stack. (Pre-condition: buffer has to be non-empty.)
2. **LEFT-ARC_r**: Add a dependency arc (w_j, r, w_i) to the arc set A , where w_i is the word second to the top of the stack and w_j is the

1. Shift $\sigma, w_i | \beta, A \rightarrow \sigma | w_i, \beta, A$
2. Left-Arc_r $\sigma | w_i | w_j, \beta, A \rightarrow \sigma | w_j, \beta, A \cup \{r(w_i, w_j)\}$
3. Right-Arc_r $\sigma | w_i | w_j, \beta, A \rightarrow \sigma | w_i, \beta, A \cup \{r(w_i, w_j)\}$

Figure 2: Transitions for Dependency Parsing.

word at the top of the stack. Remove w_i from the stack. (Pre-condition: the stack needs to contain at least two items and w_i cannot be the ROOT.)

3. **RIGHT-ARC_r**: Add a dependency arc (w_i, r, w_j) to the arc set A , where w_i is the word second to the top of the stack and w_j is the word at the top of the stack. Remove w_j from the stack. (Pre-condition: The stack needs to contain at least two items.)

A more formal definition of these three transitions is presented in Figure 2.

1.4 Neural Dependency Parsing

While there are many deep models for dependency parsing, this section focuses specifically on greedy, transition-based neural dependency parsers. This class of model has demonstrated comparable performance and significantly better efficiency than traditional feature-based discriminative dependency parsers. The primary distinction from previous models is the reliance on dense rather than sparse feature representations.

The model we will describe employs the arc-standard system for transitions, as presented in section 1.3. Ultimately, the aim of the model is to predict a transition sequence from some initial configuration c to a terminal configuration, in which the dependency parse tree is encoded. As the model is greedy, it attempts to correctly predict one transition $T \in \{\text{SHIFT}, \text{LEFT-ARC}_r, \text{RIGHT-ARC}_r\}$ at a time, based on features extracted from the current configuration $c = (\sigma, \beta, A)$. Recall, σ is the stack, β the buffer, and A the set of dependency arcs for a given sentence.

Feature Selection:

Depending on the desired complexity of the model, there is flexibility in defining the input to the neural network. The features for a given sentence S generally include some subset of:

1. S_{word} : Vector representations for some of the words in S (and their dependents) at the top of the stack σ and buffer β .
2. S_{tag} : Part-of-Speech (POS) tags for some of the words in S . POS tags comprise a small, discrete set: $\mathcal{P} = \{NN, NNP, NNS, DT, JJ, \dots\}$
3. S_{label} : The arc-labels for some of the words in S . The arc-labels comprise a small, discrete set, describing the dependency relation: $\mathcal{L} = \{amod, tmod, nsubj, csubj, dobj, \dots\}$

For each feature type, we will have a corresponding embedding matrix, mapping from the feature's one hot encoding, to a d -dimensional dense vector representation. The full embedding matrix for S_{word} is $E^w \in \mathbb{R}^{d \times N_w}$ where N_w is the dictionary/vocabulary size. Correspondingly, the POS and label embedding matrices are $E^t \in \mathbb{R}^{d \times N_t}$ and $E^l \in \mathbb{R}^{d \times N_l}$ where N_t and N_l are the number of distinct POS tags and arc labels.

Lastly, let the number of chosen elements from each set of features be denoted as n_{word} , n_{tag} , and n_{label} respectively.

Feature Selection Example:

As an example, consider the following choices for S_{word} , S_{tag} , and S_{label} .

1. S_{word} : The top 3 words on the stack and buffer: $s_1, s_2, s_3, b_1, b_2, b_3$. The first and second leftmost / rightmost children of the top two words on the stack: $lc_1(s_i), rc_1(s_i), lc_2(s_i), rc_2(s_i), i = 1, 2$. The leftmost of leftmost / rightmost of rightmost children of the top two words on the stack: $lc_1(lc_1(s_i)), rc_1(rc_1(s_i)), i = 1, 2$. In total S_{word} contains $n_w = 18$ elements.
2. S_{tag} : The corresponding POS tags for S_{tag} ($n_t = 18$).
3. S_{label} : The corresponding arc labels of words, excluding those 6 words on the stack/buffer ($n_l = 12$).

Note that we use a special NULL token for non-existent elements: when the stack and buffer are empty or dependents have not been assigned yet. For a given sentence example, we select the words, POS tags and arc labels given the schematic defined above, extract their corresponding dense feature representations produced from the embedding matrices E^w , E^t , and E^l , and concatenate these vectors into our inputs $[x^w, x^t, x^l]$. At training time we backpropagate into the dense vector representations, as well as the parameters at later layers.

Feedforward Neural Network Model:

The network contains an input layer $[x^w, x^t, x^l]$, a hidden layer, and a final softmax layer with a cross-entropy loss function. We can either define a single weight matrix in the hidden layer, to operate on a concatenation of $[x^w, x^t, x^l]$, or we can use three weight matrices $[W_1^w, W_1^t, W_1^l]$, one for each input type, as shown in Figure 3. We then apply a non-linear function and use one more affine layer $[W_2]$ so that there are an equivalent number of softmax probabilities to the number of possible transitions (the output dimension).

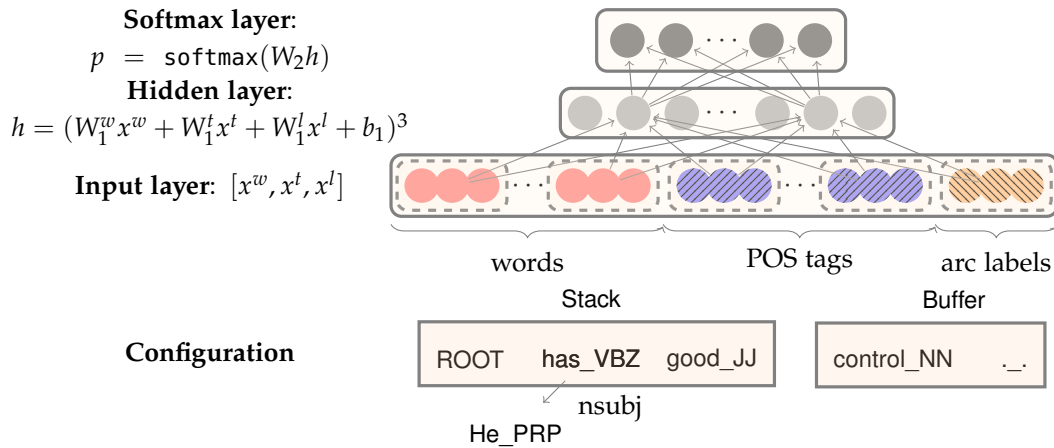


Figure 3: The neural network architecture for greedy, transition-based dependency parsing.

Note that in Figure 3, $f(x) = x^3$ is the non-linear function used.

For a more complete explanation of a greedy transition-based neural dependency parser, refer to "A Fast and Accurate Dependency Parser using Neural Networks" under Further Reading.

Further reading:

Danqi Chen, and Christopher D. Manning. "A Fast and Accurate Dependency Parser using Neural Networks." EMNLP. 2014.

Kuebler, Sandra, Ryan McDonald, and Joakim Nivre. "Dependency parsing." Synthesis Lectures on Human Language Technologies 1.1 (2009): 1-127.

CS224n: Natural Language Processing with Deep Learning ¹

Lecture Notes: Part V

Language Models, RNN, GRU and LSTM ²

Winter 2019

¹ Course Instructors: Christopher Manning, Richard Socher

² Authors: Milad Mohammadi, Rohit Mundra, Richard Socher, Lisa Wang, Amita Kamath

Keyphrases: Language Models. RNN. Bi-directional RNN. Deep RNN. GRU. LSTM.

1 Language Models

1.1 Introduction

Language models compute the probability of occurrence of a number of words in a particular sequence. The probability of a sequence of m words $\{w_1, \dots, w_m\}$ is denoted as $P(w_1, \dots, w_m)$. Since the number of words coming before a word, w_i , varies depending on its location in the input document, $P(w_1, \dots, w_m)$ is usually conditioned on a window of n previous words rather than all previous words:

$$P(w_1, \dots, w_m) = \prod_{i=1}^{i=m} P(w_i | w_1, \dots, w_{i-1}) \approx \prod_{i=1}^{i=m} P(w_i | w_{i-n}, \dots, w_{i-1}) \quad (1)$$

Equation 1 is especially useful for speech and translation systems when determining whether a word sequence is an accurate translation of an input sentence. In existing language translation systems, for each phrase / sentence translation, the software generates a number of alternative word sequences (e.g. *{I have, I had, I has, me have, me had}*) and scores them to identify the most likely translation sequence.

In machine translation, the model chooses the best word ordering for an input phrase by assigning a *goodness* score to each output word sequence alternative. To do so, the model may choose between different word ordering or word choice alternatives. It would achieve this objective by running all word sequence candidates through a probability function that assigns each a score. The sequence with the highest score is the output of the translation. For example, the machine would give a higher score to *"the cat is small"* compared to *"small the is cat"*, and a higher score to *"walking home after school"* compared to *"walking house after school"*.

1.2 n -gram Language Models

To compute the probabilities mentioned above, the count of each n -gram could be compared against the frequency of each word. This is

called an n-gram Language Model. For instance, if the model takes bi-grams, the frequency of each bi-gram, calculated via combining a word with its previous word, would be divided by the frequency of the corresponding uni-gram. Equations 2 and 3 show this relationship for bigram and trigram models.

$$p(w_2|w_1) = \frac{\text{count}(w_1, w_2)}{\text{count}(w_1)} \quad (2)$$

$$p(w_3|w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\text{count}(w_1, w_2)} \quad (3)$$

The relationship in Equation 3 focuses on making predictions based on a fixed window of context (i.e. the n previous words) used to predict the next word. But how long should the context be? In some cases, the window of past consecutive n words may not be sufficient to capture the context. For instance, consider the sentence "As the proctor started the clock, the students opened their ____". If the window only conditions on the previous three words "the students opened their", the probabilities calculated based on the corpus may suggest that the next word be "books" - however, if n had been large enough to include the "proctor" context, the probability might have suggested "exam".

This leads us to two main issues with n-gram Language Models: Sparsity and Storage.

1. Sparsity problems with n-gram Language models

Sparsity problems with these models arise due to two issues.

Firstly, note the numerator of Equation 3. If w_1 , w_2 and w_3 never appear together in the corpus, the probability of w_3 is 0. To solve this, a small δ could be added to the count for each word in the vocabulary. This is called *smoothing*.

Secondly, consider the denominator of Equation 3. If w_1 and w_2 never occurred together in the corpus, then no probability can be calculated for w_3 . To solve this, we could condition on w_2 alone. This is called *backoff*.

Increasing n makes sparsity problems worse. Typically, $n \leq 5$.

2. Storage problems with n-gram Language models

We know that we need to store the count for all n-grams we saw in the corpus. As n increases (or the corpus size increases), the model size increases as well.

1.3 Window-based Neural Language Model

The "curse of dimensionality" above was first tackled by Bengio et al in A Neural Probabilistic Language Model, which introduced the

first large-scale deep learning for natural language processing model. This model learns a *distributed representation of words*, along with the probability function for word sequences expressed in terms of these representations. Figure 1 shows the corresponding neural network architecture. The input word vectors are used by both the hidden layer and the output layer.

Equation 4 represents Figure 1 and shows the parameters of the softmax() function, consisting of the standard tanh() function (i.e. the hidden layer) as well as the linear function, $W^{(3)}x + b^{(3)}$, that captures all the previous n input word vectors.

$$\hat{y} = \text{softmax}(W^{(2)}\tanh(W^{(1)}x + b^{(1)}) + W^{(3)}x + b^{(3)}) \quad (4)$$

Note that the weight matrix $W^{(1)}$ is applied to the word vectors (solid green arrows in Figure 1), $W^{(2)}$ is applied to the hidden layer (also solid green arrow) and $W^{(3)}$ is applied to the word vectors (dashed green arrows).

A simplified version of this model can be seen in Figure 2, where the blue layer signifies concatenated word embeddings for the input words: $e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$, the red layer signifies the hidden layer: $h = f(We + b_1)$, and the green output distribution is a softmax over the vocabulary: $\hat{y} = \text{softmax}(Uh + b_2)$.

2 Recurrent Neural Networks (RNN)

Unlike the conventional translation models, where only a finite window of previous words would be considered for conditioning the language model, Recurrent Neural Networks (RNN) are capable of conditioning the model on *all* previous words in the corpus.

Figure 3 introduces the RNN architecture where each vertical rectangular box is a hidden layer at a time-step, t . Each such layer holds a number of neurons, each of which performs a linear matrix operation on its inputs followed by a non-linear operation (e.g. tanh()). At each time-step, there are two inputs to the hidden layer: the output of the previous layer h_{t-1} , and the input at that timestep x_t . The former input is multiplied by a weight matrix $W^{(hh)}$ and the latter by a weight matrix $W^{(hx)}$ to produce output features h_t , which are multiplied with a weight matrix $W^{(S)}$ and run through a softmax over the vocabulary to obtain a prediction output \hat{y} of the next word (Equations 5 and 6). The inputs and outputs of each single neuron are illustrated in Figure 4.

$$h_t = \sigma(W^{(hh)}h_{t-1} + W^{(hx)}x_{[t]}) \quad (5)$$

$$\hat{y}_t = \text{softmax}(W^{(S)}h_t) \quad (6)$$

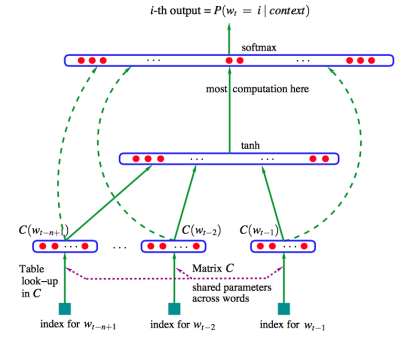


Figure 1: The first deep neural network architecture model for NLP presented by Bengio et al.

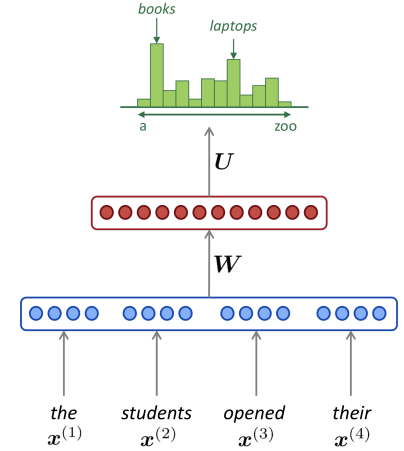


Figure 2: A simplified representation of Figure 1.

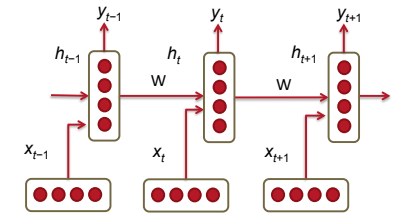


Figure 3: A Recurrent Neural Network (RNN). Three time-steps are shown.

What is interesting here is that the same weights $W^{(hh)}$ and $W^{(hx)}$ are applied repeatedly at each timestep. Thus, the number of parameters the model has to learn is less, and most importantly, is independent of the length of the input sequence - thus defeating the curse of dimensionality!

Below are the details associated with each parameter in the network:

- $x_1, \dots, x_{t-1}, x_t, x_{t+1}, \dots, x_T$: the word vectors corresponding to a corpus with T words.
- $h_t = \sigma(W^{(hh)}h_{t-1} + W^{(hx)}x_t)$: the relationship to compute the hidden layer output features at each time-step t
 - $x_t \in \mathbb{R}^d$: input word vector at time t .
 - $W^{hx} \in \mathbb{R}^{D_h \times d}$: weights matrix used to condition the input word vector, x_t
 - $W^{hh} \in \mathbb{R}^{D_h \times D_h}$: weights matrix used to condition the output of the previous time-step, h_{t-1}
 - $h_{t-1} \in \mathbb{R}^{D_h}$: output of the non-linear function at the previous time-step, $t - 1$. $h_0 \in \mathbb{R}^{D_h}$ is an initialization vector for the hidden layer at time-step $t = 0$.
 - $\sigma(\cdot)$: the non-linearity function (sigmoid here)
- $\hat{y}_t = \text{softmax}(W^{(S)}h_t)$: the output probability distribution over the vocabulary at each time-step t . Essentially, \hat{y}_t is the next predicted word given the document context score so far (i.e. h_{t-1}) and the last observed word vector $x^{(t)}$. Here, $W^{(S)} \in \mathbb{R}^{|V| \times D_h}$ and $\hat{y} \in \mathbb{R}^{|V|}$ where $|V|$ is the vocabulary.

An example of an RNN language model is shown in Figure 5. The notation in this image is slightly different: here, the equivalent of $W^{(hh)}$ is W_h , $W^{(hx)}$ is W_e , and $W^{(S)}$ is U . E converts word inputs $x^{(t)}$ to word embeddings $e^{(t)}$. The final softmax over the vocabulary shows us the probability of various options for token $x^{(5)}$, conditioned on all previous tokens. The input could be much longer than 4-5 tokens.

2.1 RNN Loss and Perplexity

The loss function used in RNNs is often the cross entropy error introduced in earlier notes. Equation 7 shows this function as the sum over the entire vocabulary at time-step t .

$$J^{(t)}(\theta) = - \sum_{j=1}^{|V|} y_{t,j} \times \log(\hat{y}_{t,j}) \quad (7)$$

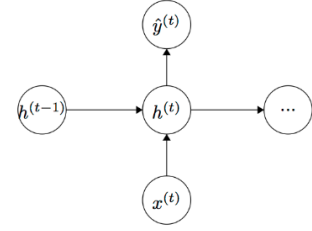


Figure 4: The inputs and outputs to a neuron of a RNN

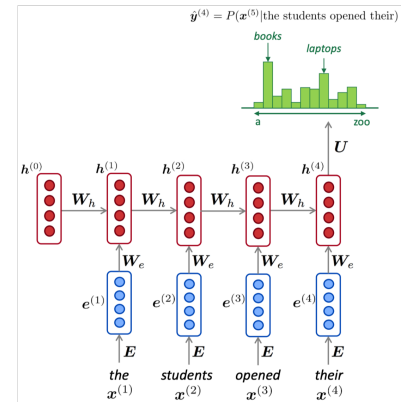


Figure 5: An RNN Language Model

The cross entropy error over a corpus of size T is:

$$J = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^{|V|} y_{t,j} \times \log(\hat{y}_{t,j}) \quad (8)$$

Equation 9 is called the *perplexity* relationship; it is basically 2 to the power of the negative log probability of the cross entropy error function shown in Equation 8. Perplexity is a measure of confusion where lower values imply more confidence in predicting the next word in the sequence (compared to the ground truth outcome).

$$\text{Perplexity} = 2^J \quad (9)$$

2.2 Advantages, Disadvantages and Applications of RNNs

RNNs have several advantages:

1. They can process input sequences of any length
2. The model size does not increase for longer input sequence lengths
3. Computation for step t can (in theory) use information from many steps back.
4. The same weights are applied to every timestep of the input, so there is symmetry in how inputs are processed

However, they also have some disadvantages:

1. Computation is slow - because it is sequential, it cannot be parallelized
2. In practice, it is difficult to access information from many steps back due to problems like vanishing and exploding gradients, which we discuss in the following subsection

The amount of memory required to run a layer of RNN is proportional to the number of words in the corpus. We can consider a sentence as a minibatch, and a sentence with k words would have k word vectors to be stored in memory. Also, the RNN must maintain two pairs of W, b matrices. As aforementioned, while the size of W could be very large, it does not scale with the size of the corpus (unlike the traditional language models). For a RNN with 1000 recurrent layers, the matrix would be 1000×1000 regardless of the corpus size.

RNNs can be used for many tasks, such as tagging (e.g. part-of-speech, named entity recognition), sentence classification (e.g. sentiment classification), and encoder modules (e.g. question answering,

machine translation, and many other tasks). In the latter two applications, we want a representation for the sentence, which we can obtain by taking the element-wise max or mean of all hidden states of the timesteps in that sentence.

Note: Figure 6 is an alternative representation of RNNs used in some publications. It represents the RNN hidden layer as a loop.

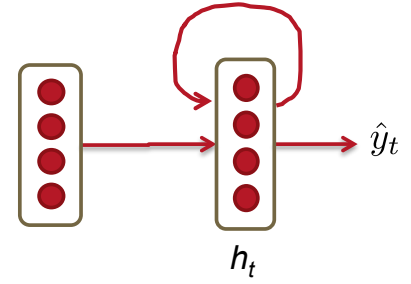


Figure 6: The illustration of a RNN as a loop over time-steps

2.3 Vanishing Gradient & Gradient Explosion Problems

Recurrent neural networks propagate weight matrices from one time-step to the next. Recall the goal of a RNN implementation is to enable propagating context information through faraway time-steps. For example, consider the following two sentences:

Sentence 1

"Jane walked into the room. John walked in too. Jane said hi to ____"

Sentence 2

"Jane walked into the room. John walked in too. It was late in the day, and everyone was walking home after a long day at work. Jane said hi to ____"

In both sentences, given their context, one can tell the answer to both blank spots is most likely "John". It is important that the RNN predicts the next word as "John", the second person who has appeared several time-steps back in both contexts. Ideally, this should be possible given what we know about RNNs so far. In practice, however, it turns out RNNs are more likely to correctly predict the blank spot in Sentence 1 than in Sentence 2. This is because during the back-propagation phase, the contribution of gradient values gradually vanishes as they propagate to earlier timesteps, as we will show below. Thus, for long sentences, the probability that "John" would be recognized as the next word reduces with the size of the context. Below, we discuss the mathematical reasoning behind the vanishing gradient problem.

Consider Equations 5 and 6 at a time-step t ; to compute the RNN error, dE/dW , we sum the error at each time-step. That is, dE_t/dW for every time-step, t , is computed and accumulated.

$$\frac{\partial E}{\partial W} = \sum_{t=1}^T \frac{\partial E_t}{\partial W} \quad (10)$$

The error for each time-step is computed through applying the chain rule differentiation to Equations 6 and 5; Equation 11 shows the corresponding differentiation. Notice dh_t/dh_k refers to the partial

derivative of h_t with respect to *all* previous k time-steps.

$$\frac{\partial E_t}{\partial W} = \sum_{k=1}^t \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W} \quad (11)$$

Equation 12 shows the relationship to compute each dh_t/dh_k ; this is simply a chain rule differentiation over all hidden layers within the $[k, t]$ time interval.

$$\frac{\partial h_t}{\partial h_k} = \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}} = \prod_{j=k+1}^t W^T \times \text{diag}[f'(h_{j-1})] \quad (12)$$

Because $h \in \mathbb{R}^{D_n}$, each $\partial h_j / \partial h_{j-1}$ is the Jacobian matrix for h :

$$\frac{\partial h_j}{\partial h_{j-1}} = \left[\frac{\partial h_j}{\partial h_{j-1,1}} \dots \frac{\partial h_j}{\partial h_{j-1,D_n}} \right] = \begin{bmatrix} \frac{\partial h_{j,1}}{\partial h_{j-1,1}} & \cdot & \cdot & \cdot & \frac{\partial h_{j,1}}{\partial h_{j-1,D_n}} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \frac{\partial h_{j,D_n}}{\partial h_{j-1,1}} & \cdot & \cdot & \cdot & \frac{\partial h_{j,D_n}}{\partial h_{j-1,D_n}} \end{bmatrix} \quad (13)$$

Putting Equations 10, 11, 12 together, we have the following relationship.

$$\frac{\partial E}{\partial W} = \sum_{t=1}^T \sum_{k=1}^t \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \left(\prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}} \right) \frac{\partial h_k}{\partial W} \quad (14)$$

Equation 15 shows the norm of the Jacobian matrix relationship in Equation 13. Here, β_W and β_h represent the upper bound values for the two matrix norms. The norm of the partial gradient at each time-step, t , is therefore, calculated through the relationship shown in Equation 15.

$$\left\| \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq \|W^T\| \|\text{diag}[f'(h_{j-1})]\| \leq \beta_W \beta_h \quad (15)$$

The norm of both matrices is calculated through taking their L2-norm. The norm of $f'(h_{j-1})$ can only be as large as 1 given the sigmoid non-linearity function.

$$\left\| \frac{\partial h_t}{\partial h_k} \right\| = \left\| \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq (\beta_W \beta_h)^{t-k} \quad (16)$$

The exponential term $(\beta_W \beta_h)^{t-k}$ can easily become a very small or large number when $\beta_W \beta_h$ is much smaller or larger than 1 and $t - k$ is sufficiently large. Recall a large $t - k$ evaluates the cross entropy error due to faraway words. The contribution of faraway words to predicting the next word at time-step t diminishes when the gradient vanishes early on.

During experimentation, once the gradient value grows extremely large, it causes an overflow (i.e. NaN) which is easily detectable at runtime; this issue is called the *Gradient Explosion Problem*. When the gradient value goes to zero, however, it can go undetected while drastically reducing the learning quality of the model for far-away words in the corpus; this issue is called the *Vanishing Gradient Problem*. Due to vanishing gradients, we don't know whether there is no dependency between steps t and $t + n$ in the data, or we just cannot capture the true dependency due to this issue.

To gain practical intuition about the vanishing gradient problem, you may visit the following [example website](#).

2.4 Solution to the Exploding & Vanishing Gradients

Now that we gained intuition about the nature of the vanishing gradients problem and how it manifests itself in deep neural networks, let us focus on a simple and practical heuristic to solve these problems.

To solve the problem of exploding gradients, Thomas Mikolov first introduced a simple heuristic solution that *clips* gradients to a small number whenever they explode. That is, whenever they reach a certain threshold, they are set back to a small number as shown in Algorithm 1.

```

 $\hat{g} \leftarrow \frac{\partial E}{\partial W}$ 
if  $\|\hat{g}\| \geq threshold$  then
     $\hat{g} \leftarrow \frac{threshold}{\|\hat{g}\|} \hat{g}$ 
end if
    
```

Algorithm 1: Pseudo-code for norm clipping in the gradients whenever they explode

Figure 7 visualizes the effect of gradient clipping. It shows the decision surface of a small recurrent neural network with respect to its W matrix and its bias terms, b . The model consists of a single unit of recurrent neural network running through a small number of time-steps; the solid arrows illustrate the training progress on each gradient descent step. When the gradient descent model hits the high error wall in the objective function, the gradient is pushed off to a far-away location on the decision surface. The clipping model produces the dashed line where it instead pulls back the error gradient to somewhere close to the original gradient landscape.

To solve the problem of vanishing gradients, we introduce two techniques. The first technique is that instead of initializing $W^{(hh)}$ randomly, start off from an identity matrix initialization.

The second technique is to use the Rectified Linear Units (ReLU) instead of the sigmoid function. The derivative for the ReLU is either 0 or

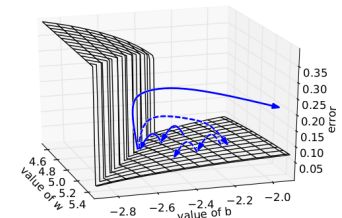


Figure 7: Gradient explosion clipping visualization

1. This way, gradients would flow through the neurons whose derivative is 1 without getting attenuated while propagating back through time-steps.

2.5 Deep Bidirectional RNNs

So far, we have focused on RNNs that condition on past words to predict the next word in the sequence. It is possible to make predictions based on future words by having the RNN model read through the corpus backwards. Irsoy et al. shows a bi-directional deep neural network; at each time-step, t , this network maintains two hidden layers, one for the left-to-right propagation and another for the right-to-left propagation. To maintain two hidden layers at any time, this network consumes twice as much memory space for its weight and bias parameters. The final classification result, \hat{y}_t , is generated through combining the score results produced by both RNN hidden layers. Figure 8 shows the bi-directional network architecture, and Equations 17 and 18 show the mathematical formulation behind setting up the bi-directional RNN hidden layer. The only difference between these two relationships is in the direction of recursing through the corpus. Equation 19 shows the classification relationship used for predicting the next word via summarizing past and future word representations.

$$\vec{h}_t = f(\vec{W}x_t + \vec{V}\vec{h}_{t-1} + \vec{b}) \quad (17)$$

$$\overleftarrow{h}_t = f(\overleftarrow{W}x_t + \overleftarrow{V}\overleftarrow{h}_{t+1} + \overleftarrow{b}) \quad (18)$$

$$\hat{y}_t = g(Uh_t + c) = g(U[\vec{h}_t; \overleftarrow{h}_t] + c) \quad (19)$$

RNNs can also be multi-layered. Figure 9 shows a multi-layer bi-directional RNN where each lower layer feeds the next layer. As shown in this figure, in this network architecture, at time-step t each intermediate neuron receives one set of parameters from the previous time-step (in the same RNN layer), and two sets of parameters from the previous RNN hidden layer; one input comes from the left-to-right RNN and the other from the right-to-left RNN.

To construct a Deep RNN with L layers, the above relationships are modified to the relationships in Equations 20 and 21 where the input to each intermediate neuron at level i is the output of the RNN at layer $i - 1$ at the same time-step, t . The output, \hat{y} , at each time-step is the result of propagating input parameters through all hidden layers (Equation 22).

$$\vec{h}_t^{(i)} = f(\vec{W}^{(i)}h_t^{(i-1)} + \vec{V}^{(i)}\vec{h}_{t-1}^{(i)} + \vec{b}^{(i)}) \quad (20)$$

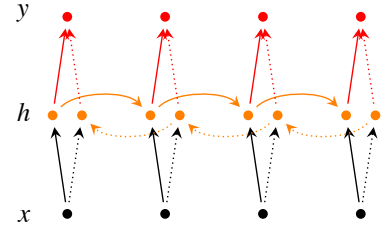


Figure 8: A bi-directional RNN model

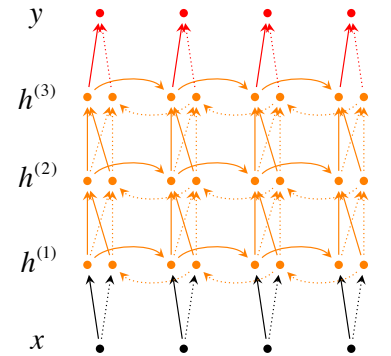


Figure 9: A deep bi-directional RNN with three RNN layers.

$$\overleftarrow{h}_t^{(i)} = f(\overleftarrow{W}^{(i)} h_t^{(i-1)} + \overleftarrow{V}^{(i)} \overleftarrow{h}_{t+1}^{(i)} + \overleftarrow{b}^{(i)}) \quad (21)$$

$$\hat{y}_t = g(Uh_t + c) = g(U[\overrightarrow{h}_t^{(L)}; \overleftarrow{h}_t^{(L)}] + c) \quad (22)$$

2.6 Application: RNN Translation Model

Traditional translation models are quite complex; they consist of numerous machine learning algorithms applied to different stages of the language translation pipeline. In this section, we discuss the potential for adopting RNNs as a replacement to traditional translation modules. Consider the RNN example model shown in Figure 10; here, the German phrase *Echt dicke Kiste* is translated to *Awesome sauce*. The first three hidden layer time-steps *encode* the German language words into some language word features (h_3). The last two time-steps *decode* h_3 into English word outputs. Equation 23 shows the relationship for the Encoder stage and Equations 24 and 25 show the equation for the Decoder stage.

$$h_t = \phi(h_{t-1}, x_t) = f(W^{(hh)}h_{t-1} + W^{(hx)}x_t) \quad (23)$$

$$h_t = \phi(h_{t-1}) = f(W^{(hh)}h_{t-1}) \quad (24)$$

$$y_t = \text{softmax}(W^{(S)}h_t) \quad (25)$$

One may naively assume this RNN model along with the cross-entropy function shown in Equation 26 can produce high-accuracy translation results. In practice, however, several extensions are to be added to the model to improve its translation accuracy performance.

$$\max_{\theta} \frac{1}{N} \sum_{n=1}^N \log p_{\theta}(y^{(n)} | x^{(n)}) \quad (26)$$

Extension I: train different RNN weights for encoding and decoding. This decouples the two units and allows for more accuracy prediction of each of the two RNN modules. This means the $\phi(\cdot)$ functions in Equations 23 and 24 would have different $W^{(hh)}$ matrices.

Extension II: compute every hidden state in the decoder using three different inputs:

- The previous hidden state (standard)
- Last hidden layer of the encoder ($c = h_T$ in Figure 11)

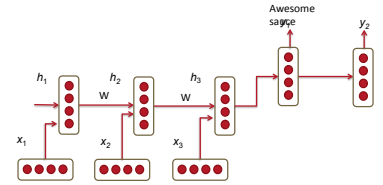


Figure 10: A RNN-based translation model. The first three RNN hidden layers belong to the source language model encoder, and the last two belong to the destination language model decoder.

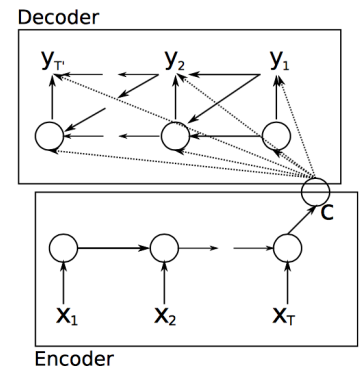


Figure 11: Language model with three inputs to each decoder neuron: (h_{t-1}, c, y_{t-1})

- Previous predicted output word, \hat{y}_{t-1}

Combining the above three inputs transforms the ϕ function in the decoder function of Equation 24 to the one in Equation 27. Figure 11 illustrates this model.

$$h_t = \phi(h_{t-1}, c, y_{t-1}) \quad (27)$$

Extension III: train deep recurrent neural networks using multiple RNN layers as discussed earlier in this chapter. Deeper layers often improve prediction accuracy due to their higher learning capacity. Of course, this implies a large training corpus must be used to train the model.

Extension IV: train bi-directional encoders to improve accuracy similar to what was discussed earlier in this chapter.

Extension V: given a word sequence $A B C$ in German whose translation is $X Y$ in English, instead of training the RNN using $A B C \rightarrow X Y$, train it using $C B A \rightarrow X Y$. The intuition behind this technique is that A is more likely to be translated to X . Thus, given the vanishing gradient problem discussed earlier, reversing the order of the input words can help reduce the error rate in generating the output phrase.

3 Gated Recurrent Units

Beyond the extensions discussed so far, RNNs have been found to perform better with the use of more complex units for activation. So far, we have discussed methods that transition from hidden state h_{t-1} to h_t using an affine transformation and a point-wise nonlinearity. Here, we discuss the use of a gated activation function thereby modifying the RNN architecture. What motivates this? Well, although RNNs can theoretically capture long-term dependencies, they are very hard to actually train to do this. Gated recurrent units are designed in a manner to have more persistent memory thereby making it easier for RNNs to capture long-term dependencies. Let us see mathematically how a GRU uses h_{t-1} and x_t to generate the next hidden state h_t . We will then dive into the intuition of this architecture.

$$\begin{aligned} z_t &= \sigma(W^{(z)}x_t + U^{(z)}h_{t-1}) && \text{(Update gate)} \\ r_t &= \sigma(W^{(r)}x_t + U^{(r)}h_{t-1}) && \text{(Reset gate)} \\ \tilde{h}_t &= \tanh(r_t \circ U h_{t-1} + W x_t) && \text{(New memory)} \\ h_t &= (1 - z_t) \circ \tilde{h}_t + z_t \circ h_{t-1} && \text{(Hidden state)} \end{aligned}$$

The above equations can be thought of a GRU's four fundamental operational stages and they have intuitive interpretations that make this model much more intellectually satisfying (see Figure 12):

1. **New memory generation:** A new memory \tilde{h}_t is the consolidation of a new input word x_t with the past hidden state h_{t-1} . Anthropomorphically, this stage is the one who knows the recipe of combining a newly observed word with the past hidden state h_{t-1} to summarize this new word in light of the contextual past as the vector \tilde{h}_t .
2. **Reset Gate:** The reset signal r_t is responsible for determining how important h_{t-1} is to the summarization \tilde{h}_t . The reset gate has the ability to completely diminish past hidden state if it finds that h_{t-1} is irrelevant to the computation of the new memory.
3. **Update Gate:** The update signal z_t is responsible for determining how much of h_{t-1} should be carried forward to the next state. For instance, if $z_t \approx 1$, then h_{t-1} is almost entirely copied out to h_t . Conversely, if $z_t \approx 0$, then mostly the new memory \tilde{h}_t is forwarded to the next hidden state.
4. **Hidden state:** The hidden state h_t is finally generated using the past hidden input h_{t-1} and the new memory generated \tilde{h}_t with the advice of the update gate.

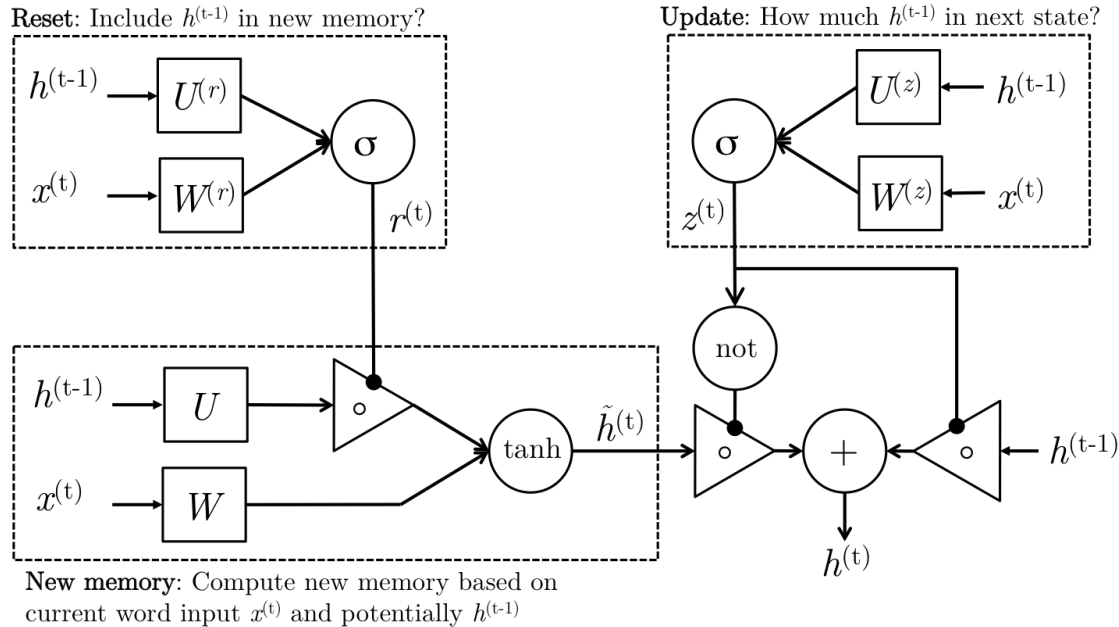


Figure 12: The detailed internals of a GRU

It is important to note that to train a GRU, we need to learn all the different parameters: $W, U, W^{(r)}, U^{(r)}, W^{(z)}, U^{(z)}$. These follow the same backpropagation procedure we have seen in the past.

4 Long-Short-Term-Memories

Long-Short-Term-Memories are another type of complex activation unit that differ a little from GRUs. The motivation for using these is similar to those for GRUs however the architecture of such units does differ. Let us first take a look at the mathematical formulation of LSTM units before diving into the intuition behind this design:

$$\begin{aligned}
 i_t &= \sigma(W^{(i)}x_t + U^{(i)}h_{t-1}) && \text{(Input gate)} \\
 f_t &= \sigma(W^{(f)}x_t + U^{(f)}h_{t-1}) && \text{(Forget gate)} \\
 o_t &= \sigma(W^{(o)}x_t + U^{(o)}h_{t-1}) && \text{(Output/Exposure gate)} \\
 \tilde{c}_t &= \tanh(W^{(c)}x_t + U^{(c)}h_{t-1}) && \text{(New memory cell)} \\
 c_t &= f_t \circ c_{t-1} + i_t \circ \tilde{c}_t && \text{(Final memory cell)} \\
 h_t &= o_t \circ \tanh(c_t) &&
 \end{aligned}$$

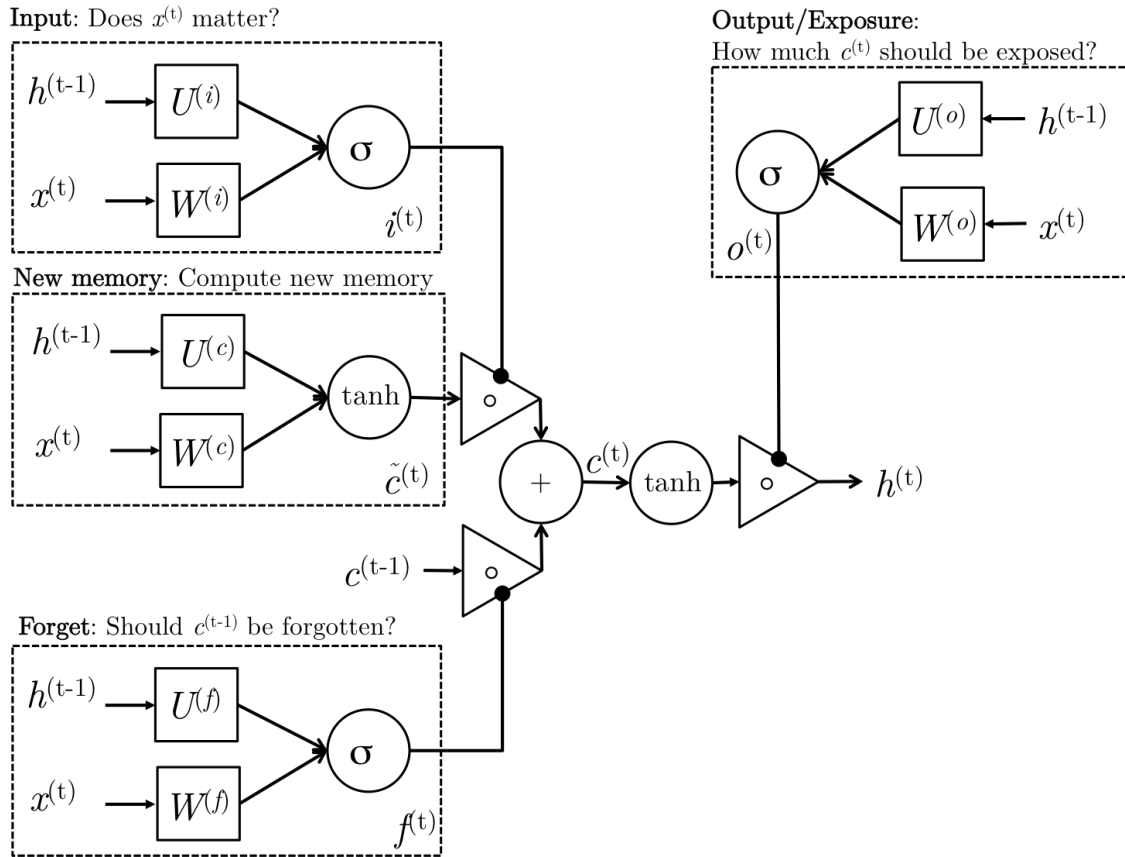


Figure 13: The detailed internals of a LSTM

We can gain intuition of the structure of an LSTM by thinking of its architecture as the following stages:

1. **New memory generation:** This stage is analogous to the new memory generation stage we saw in GRUs. We essentially use the input word x_t and the past hidden state h_{t-1} to generate a new memory \tilde{c}_t which includes aspects of the new word $x^{(t)}$.
2. **Input Gate:** We see that the new memory generation stage doesn't check if the new word is even important before generating the new memory – this is exactly the input gate's function. The input gate uses the input word and the past hidden state to determine whether or not the input is worth preserving and thus is used to gate the new memory. It thus produces i_t as an indicator of this information.
3. **Forget Gate:** This gate is similar to the input gate except that it does not make a determination of usefulness of the input word – instead it makes an assessment on whether the past memory cell is useful for the computation of the current memory cell. Thus, the forget gate looks at the input word and the past hidden state and produces f_t .
4. **Final memory generation:** This stage first takes the advice of the forget gate f_t and accordingly forgets the past memory c_{t-1} . Similarly, it takes the advice of the input gate i_t and accordingly gates the new memory \tilde{c}_t . It then sums these two results to produce the final memory c_t .
5. **Output/Exposure Gate:** This is a gate that does not explicitly exist in GRUs. It's purpose is to separate the final memory from the hidden state. The final memory c_t contains a lot of information that is not necessarily required to be saved in the hidden state. Hidden states are used in every single gate of an LSTM and thus, this gate makes the assessment regarding what parts of the memory c_t needs to be exposed/present in the hidden state h_t . The signal it produces to indicate this is o_t and this is used to gate the point-wise tanh of the memory.

CS224n: Natural Language Processing with Deep Learning¹

Lecture Notes: Part VI

Neural Machine Translation, Seq2seq and Attention²

Winter 2019

¹ Course Instructors: Christopher Manning, Richard Socher

² Authors: Guillaume Genthial, Lucas Liu, Barak Oshri, Kushal Ranjan

Keyphrases: Seq2Seq and Attention Mechanisms, Neural Machine Translation, Speech Processing

1 Neural Machine Translation with Seq2Seq

So far in this class, we've dealt with problems of predicting a single output: an NER label for a word, the single most likely next word in a sentence given the past few, and so on. However, there's a whole class of NLP tasks that rely on *sequential output*, or outputs that are sequences of potentially varying length. For example,

- **Translation:** taking a sentence in one language as input and outputting the same sentence in another language.
- **Conversation:** taking a statement or question as input and responding to it.
- **Summarization:** taking a large body of text as input and outputting a summary of it.

In these notes, we'll look at sequence-to-sequence models, a deep learning-based framework for handling these types of problems. This framework proved to be very effective, and has, in fewer than 3 years, become the standard for machine translation.

1.1 Brief Note on Historical Approaches

In the past, translation systems were based on probabilistic models constructed from:

- a **translation model**, telling us what a sentence/phrase in a source language most likely translates into
- a **language model**, telling us how likely a given sentence/phrase is overall.

These components were used to build translation systems based on words or phrases. As you might expect, a naive word-based system would completely fail to capture differences in ordering between languages (e.g. where negation words go, location of subject vs. verb in a sentence, etc).

Phrase-based systems were most common prior to Seq2Seq. A phrase-based translation system can consider inputs and outputs in terms of sequences of phrases and can handle more complex syntaxes than word-based systems. However, long-term dependencies are still difficult to capture in phrase-based systems.

The advantage that Seq2Seq brought to the table, especially with its use of LSTMs, is that modern translation systems can generate arbitrary output sequences after seeing the *entire* input. They can even focus in on specific parts of the input automatically to help generate a useful translation.

1.2 Sequence-to-sequence Basics

Sequence-to-sequence, or "Seq2Seq", is a relatively new paradigm, with its first published usage in 2014 for English-French translation³. At a high level, a sequence-to-sequence model is an end-to-end model made up of two recurrent neural networks:

- an *encoder*, which takes the model's input sequence as input and encodes it into a fixed-size "context vector", and
- a *decoder*, which uses the context vector from above as a "seed" from which to generate an output sequence.

For this reason, Seq2Seq models are often referred to as "encoder-decoder models." We'll look at the details of these two networks separately.

³ Sutskever et al. 2014, "Sequence to Sequence Learning with Neural Networks"

1.3 Seq2Seq architecture - encoder

The encoder network's job is to read the input sequence to our Seq2Seq model and generate a fixed-dimensional context vector C for the sequence. To do so, the encoder will use a recurrent neural network cell – usually an LSTM – to read the input tokens one at a time. The final hidden state of the cell will then become C . However, because it's so difficult to compress an arbitrary-length sequence into a single fixed-size vector (especially for difficult tasks like translation), the encoder will usually consist of *stacked* LSTMs: a series of LSTM "layers" where each layer's outputs are the input sequence to the next layer. The *final* layer's LSTM hidden state will be used as C .

Seq2Seq encoders will often do something strange: they will process the input sequence *in reverse*. This is actually done on purpose. The idea is that, by doing this, the *last* thing that the encoder sees will (roughly) corresponds to the *first* thing that the model outputs; this makes it easier for the decoder to "get started" on the output, which makes then gives the decoder an easier time generating a

proper output sentence. In the context of translation, we're allowing the network to translate the first few words of the input as soon as it sees them; once it has the first few words translated correctly, it's much easier to go on to construct a correct sentence than it is to do so from scratch. See Fig. 1 for an example of what such an encoder network might look like.

1.4 Seq2Seq architecture - decoder

The decoder is also an LSTM network, but its usage is a little more complex than the encoder network. Essentially, we'd like to use it as a language model that's "aware" of the words that it's generated so far *and* of the input. To that end, we'll keep the "stacked" LSTM architecture from the encoder, but we'll initialize the hidden state of our first layer with the context vector from above; the decoder will literally use the context of the input to generate an output.

Once the decoder is set up with its context, we'll pass in a special token to signify the start of output generation; in literature, this is usually an $\langle \text{EOS} \rangle$ token appended to the end of the input (there's also one at the end of the output). Then, we'll run all three layers of LSTM, one after the other, following up with a softmax on the final layer's output to generate the first output word. Then, we *pass that word into the first layer*, and repeat the generation. This is how we get the LSTMs to act like a language model. See Fig. 2 for an example of a decoder network.

Once we have the output sequence, we use the same learning strategy as usual. We define a loss, the cross entropy on the prediction sequence, and we minimize it with a gradient descent algorithm and back-propagation. Both the encoder and decoder are trained at the same time, so that they both learn the same context vector representation.

1.5 Recap & Basic NMT Example

Note that there is no connection between the lengths of the input and output; any length input can be passed in and any length output can be generated. However, Seq2Seq models are known to lose effectiveness on very long inputs, a consequence of the practical limits of LSTMs.

To recap, let's think about what a Seq2Seq model does in order to translate the English "what is your name" into the French "comment t'appelles tu". First, we start with 4 one-hot vectors for the input. These inputs may or may not (for translation, they usually are) embedded into a dense vector representation. Then, a stacked LSTM network reads the sequence in reverse and *encodes* it into a context

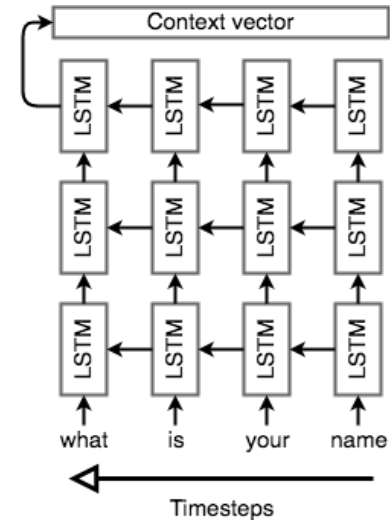


Figure 1: Example of a Seq2Seq encoder network. This model may be used to translate the English sentence "what is your name?" Note that the input tokens are read in reverse. Note that the network is unrolled; each column is a timestep and each row is a single layer, so that horizontal arrows correspond to hidden states and vertical arrows are LSTM inputs/outputs.

vector. This context vector is a vector space representation of the notion of asking someone for their name. It's used to initialize the first layer of another stacked LSTM. We run one step of each layer of this network, perform softmax on the last layer's output, and use that to select our first output word. This word is fed back into the network as input, and the rest of the sentence "comment t'appelles tu" is *de-coded* in this fashion. During backpropagation, the encoder's LSTM weights are updated so that it learns a better vector space representation for sentences, while the decoder's LSTM weights are trained to allow it to generate grammatically correct sentences that are relevant to the context vector.

1.6 Bidirectional RNNs

Recall from earlier in this class that dependencies in sentences don't just work in one direction; a word can have a dependency on another word before *or* after it. The formulation of Seq2Seq that we've talked about so far doesn't account for that; at any timestep, we're only considering information (via the LSTM hidden state) from words *before* the current word. For NMT, we need to be able to effectively encode any input, regardless of dependency directions within that input, so this won't cut it.

Bidirectional RNNs fix this problem by traversing a sequence in both directions and concatenating the resulting outputs (both cell outputs and final hidden states). For every RNN cell, we simply add another cell but feed inputs to it in the opposite direction; the output o_t corresponding to the t 'th word is the concatenated vector $[o_t^{(f)} \ o_t^{(b)}]$, where $o_t^{(f)}$ is the output of the forward-direction RNN on word t and $o_t^{(b)}$ is the corresponding output from the reverse-direction RNN. Similarly, the final hidden state is $h = [h^{(f)} \ h^{(b)}]$, where $h^{(f)}$ is the final hidden state of the forward RNN and $h^{(b)}$ is the final hidden state of the reverse RNN. See Fig. 6 for an example of a bidirectional LSTM encoder.

2 Attention Mechanism

2.1 Motivation

When you hear the sentence "the ball is on the field," you don't assign the same importance to all 6 words. You primarily take note of the words "ball," "on," and "field," since those are the words that are most "important" to you. Similarly, Bahdanau et al. noticed the flaw in using the final RNN hidden state as the single "context vector" for sequence-to-sequence models: often, different parts of an input have

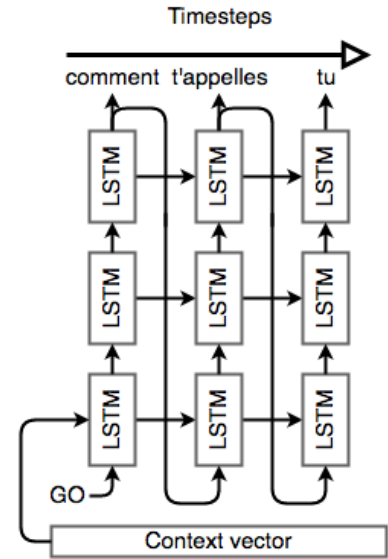


Figure 2: Example of a Seq2Seq decoder network. This decoder is decoding the context vector for "what is your name" (see Fig. 1 into its French translation, "comment t'appelles tu?"). Note the special "GO" token used at the start of generation, and that generation is in the forward direction as opposed to the input which is read in reverse. Note also that the input and output do not need to be the same length.

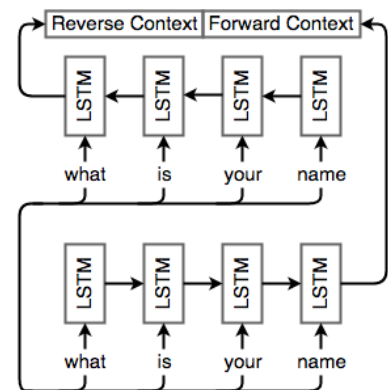


Figure 3: Example of a single-layer bidirectional LSTM encoder network. Note that the input is fed into two different LSTM layers, but in different directions, and the hidden states are concatenated to get the final context vector.

different levels of significance. Moreover, different parts of the output may even consider different parts of the input "important." For example, in translation, the first word of output is *usually* based on the first few words of the input, but the last word is likely based on the last few words of input.

Attention mechanisms make use of this observation by providing the decoder network with a look at the *entire input sequence* at every decoding step; the decoder can then decide what input words are important at any point in time. There are many types of encoder mechanisms, but we'll examine the one introduced by Bahdanau et al.⁴,

⁴ Bahdanau et al. 2014, "Neural Machine Translation by Jointly Learning to Align and Translate"

2.2 Bahdanau et al. NMT model

Remember that our seq2seq model is made of two parts, an **encoder** that encodes the input sentence, and a **decoder** that leverages the information extracted by the decoder to produce the translated sentence. Basically, our input is a sequence of words x_1, \dots, x_n that we want to translate, and our target sentence is a sequence of words y_1, \dots, y_m .

1. Encoder

Let (h_1, \dots, h_n) be the hidden vectors representing the input sentence. These vectors are the output of a bi-LSTM for instance, and capture contextual representation of each word in the sentence.

2. Decoder

We want to compute the hidden states s_i of the decoder using a recursive formula of the form

$$s_i = f(s_{i-1}, y_{i-1}, c_i)$$

where s_{i-1} is the previous hidden vector, y_{i-1} is the generated word at the previous step, and c_i is a context vector that capture the context from the original sentence that is relevant to the time step i of the decoder.

The context vector c_i captures relevant information for the i -th decoding time step (unlike the standard Seq2Seq in which there's only one context vector). For each hidden vector from the original sentence h_j , compute a score

$$e_{i,j} = a(s_{i-1}, h_j)$$

where a is any function with values in \mathbb{R} , for instance a single layer fully-connected neural network. Then, we end up with a

sequence of scalar values $e_{i,1}, \dots, e_{i,n}$. Normalize these scores into a vector $\alpha_i = (\alpha_{i,1}, \dots, \alpha_{i,n})$, using a *softmax* layer.

$$\alpha_{i,j} = \frac{\exp(e_{i,j})}{\sum_{k=1}^n \exp(e_{i,k})}$$

Then, compute the context vector c_i as the weighted average of the hidden vectors from the original sentence

$$c_i = \sum_{j=1}^n \alpha_{i,j} h_j$$

Intuitively, this vector captures the relevant contextual information from the original sentence for the i -th step of the decoder.

The vector α_i is called the *attention* vector

The context vector is extracted thanks to the attention vector and captures the relevant context

2.3 Connection with translation alignment

The attention-based model learns to assign significance to different parts of the input for each step of the output. In the context of translation, attention can be thought of as "alignment." Bahdanau et al. argue that the attention scores α_{ij} at decoding step i signify the words in the source sentence that align with word i in the target. Noting this, we can use attention scores to build an alignment table – a table mapping words in the source to corresponding words in the target sentence – based on the learned encoder and decoder from our Seq2Seq NMT system.

	Comment	t'	appelles	tu	<end>
What					
is					
your					
name					
<end>					

Figure 4: Example of an alignment table

2.4 Performance on long sentences

The major advantage of attention-based models is their ability to efficiently translate long sentences. As the size of the input grows, models that do not use attention will miss information and precision if they only use the final representation. Attention is a clever way to fix this issue and experiments indeed confirm the intuition.

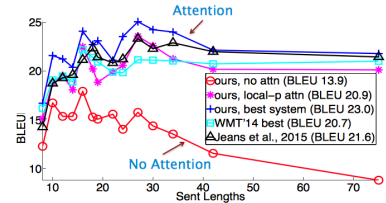


Figure 5: Performance on long sentence of different NMT models - image taken from Luong et al.

3 Other Models

3.1 Huong et al. NMT model

We present a variant of this first model, with two different mechanisms of attention, from Luong et al.⁵.

- **Global attention** We run our vanilla Seq2Seq NMT. We call the hidden states given by the encoder h_1, \dots, h_n , and the hidden states of the decoder $\tilde{h}_1, \dots, \tilde{h}_n$. Now, for each \tilde{h}_i , we compute an

⁵ *Effective Approaches to Attention-based Neural Machine Translation* by Minh-Thang Luong, Hieu Pham and Christopher D. Manning

attention vector over the encoder hidden. We can use one of the following scoring functions:

$$\text{score}(h_i, \tilde{h}_j) = \begin{cases} h_i^T \tilde{h}_j \\ h_i^T W \tilde{h}_j \\ W[h_i, \tilde{h}_j] \end{cases} \in \mathbb{R}$$

Now that we have a vector of scores, we can compute a context vector in the same way as Bahdanau et al. First, we normalize the scores via a softmax layer to obtain a vector $\alpha_i = (\alpha_{i,1}, \dots, \alpha_{i,n})$,

where $\alpha_{i,j} = \frac{\exp(\text{score}(h_i, \tilde{h}_j))}{\sum_{k=1}^n \exp(\text{score}(h_i, \tilde{h}_k))}$

$$c_i = \sum_{j=1}^n \alpha_{i,j} h_j$$

and we can use the context vector and the hidden state to compute a new vector for the i -th time step of the decoder

$$\tilde{h}_i = f([\tilde{h}_i, c_i])$$

The final step is to use the \tilde{h}_i to make the final prediction of the decoder. To address the issue of coverage, Luong et al. also use an input-feeding approach. The attentional vectors \tilde{h}_i are fed as input to the decoder, instead of the final prediction. This is similar to Bahdanau et al., who use the context vectors to compute the hidden vectors of the decoder.

- **Local attention** the model predicts an aligned position in the input sequence. Then, it computes a context vector using a window centered on this position. The computational cost of this attention step is constant and does not explode with the length of the sentence.

The main takeaway of this discussion is to show that there are lots of ways of doing attention.

3.2 Google's new NMT

As a brief aside, Google recently made a major breakthrough for NMT via their own translation system⁶. Rather than maintain a full Seq2Seq model for every pair of language that they support – each of which would have to be trained individually, which is a tremendous feat in terms of both data and compute time required – they built a single system that can translate between any two languages. This is a Seq2Seq model that accepts as input a sequence of words *and* a token

⁶ Johnson et al. 2016, "Google's Multilingual Neural Machine Translation System: Enabling Zero-Shot Translation"

specifying what language to translate into. The model uses shared parameters to translate into any target language.

The new multilingual model not only improved their translation performance, it also enabled "zero-shot translation," in which we can translate between two languages *for which we have no translation training data*. For instance, if we only had examples of Japanese-English translations and Korean-English translations, Google's team found that the multilingual NMT system trained on this data could actually generate reasonable Japanese-Korean translations. The powerful implication of this finding is that part of the decoding process is not language-specific, and the model is in fact maintaining an internal representation of the input/output sentences independent of the actual languages involved.

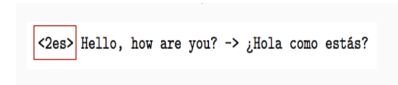


Figure 6: Example of Google's system

3.3 More advanced papers using attention

- *Show, Attend and Tell: Neural Image Caption Generation with Visual Attention* by Kelvin Xu, Jimmy Lei Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard S. Zemel and Yoshua Bengio. This paper learns words/image alignment.
- *Modeling Coverage for Neural Machine Translation* by Zhaopeng Tu, Zhengdong Lu, Yang Liu, Xiaohua Liu and Hang Li. Their model uses a coverage vector that takes into account the attention history to help future attention.
- *Incorporating Structural Alignment Biases into an Attentional Neural Translation Model* by Cohn, Hoang, Vymolova, Yao, Dyer, Haffari. This paper improves the attention by incorporating other traditional linguistic ideas.

4 Sequence model decoders

Another approach to machine translation comes from statistical machine translation. Consider a model that computes the probability $\mathbb{P}(\bar{s}|s)$ of a translation \bar{s} given the original sentence s . We want to pick the translation \bar{s}^* that has the best probability. In other words, we want

$$\bar{s}^* = \operatorname{argmax}_{\bar{s}} (\mathbb{P}(\bar{s}|s))$$

As the search space can be huge, we need to shrink its size. Here is a list of sequence model decoders (both good ones and bad ones).

- **Exhaustive search** : this is the simplest idea. We compute the probability of every possible sequence, and we chose the sequence

with the highest probability. However, this technique does not scale *at all* to large outputs as the search space is exponential in the size of the input. Decoding in this case is an NP-complete problem.

- **Ancestral sampling** : at time step t , we sample x_t based on the conditional probability of the word at step t given the past. In other words,

$$x_t \sim \mathbb{P}(x_t | x_1, \dots, x_n)$$

Theoretically, this technique is efficient and asymptotically exact. However, in practice, it can have low performance and high variance.

- **Greedy Search** : At each time step, we pick the most probable token. In other words

$$x_t = \operatorname{argmax}_{\tilde{x}_t} \mathbb{P}(\tilde{x}_t | x_1, \dots, x_n)$$

This technique is efficient and natural, however it explores a small part of the search space and if we make a mistake at one time step, the rest of the sentence could be heavily impacted.

- **Beam search** : the idea is to maintain K candidates at each time step.

$$\mathcal{H}_t = \{(x_1^1, \dots, x_t^1), \dots, (x_1^K, \dots, x_t^K)\}$$

and compute \mathcal{H}_{t+1} by expanding \mathcal{H}_t and keeping the best K candidates. In other words, we pick the best K sequence in the following set

$$\tilde{\mathcal{H}}_{t+1} = \bigcup_{k=1}^K \mathcal{H}_{t+1}^k$$

where

$$\mathcal{H}_{t+1}^k = \{(x_1^k, \dots, x_t^k, v_1), \dots, (x_1^k, \dots, x_t^k, v_{|V|})\}$$

As we increase K , we gain precision and we are asymptotically exact. However, the improvement is not monotonic and we can set a K that combines reasonable performance and computational efficiency. For this reason, beam search is the most commonly used technique in NMT.

5 *Evaluation of Machine Translation Systems*

Now that we know the basics about machine translation systems, we discuss some ways that these models are evaluated. Evaluating the quality of translations is a notoriously tricky and subjective task. In real-life, if you give a paragraph of text to ten different translators, you will get back ten different translations. Translations are imperfect and noisy in practice. They attend to different information and emphasize different meanings. One translation can preserve metaphors and the integrity of long-ranging ideas, while the other can achieve a more faithful reconstruction of syntax and style, attempting a word-to-word translation. Note that this flexibility is not a burden; it is a testament to the complexity of language and our abilities to decode and interpret meaning, and is a wonderful aspect of our communicative faculty.

At this point, you should note that there is a difference between the objective *loss* function of your model and the *evaluation* methods we are going to discuss. Since loss functions are in essence an evaluation of your model prediction, it can be easy to confuse the two ideas. The evaluation metrics ahead offer a final, summative assessment of your model against some measurement criterion, and no one measurement is superior to all others, though some have clear advantages and majority preference.

Evaluating the quality of machine learning translations has become its own entire research area, with many proposals like TER, METEOR, MaxSim, SEPIA, and RTE-MT. We will focus in these notes on two baseline evaluation methods and BLEU.

5.1 *Human Evaluation*

The first and maybe least surprising method is to have people manually evaluate the correctness, adequacy, and fluency of your system. Like the Turing Test, if you can fool a human into not being able to distinguish a human-made translation with your system translation, your model passes the test for looking like a real-life sentence! The obvious problem with this method is that it is costly and inefficient, though it remains the gold standard for machine translation.

5.2 *Evaluation against another task*

A common way of evaluating machine learning models that output a useful *representation* of some data (a representation being a translation or summary) is that if your predictions are useful for solving some challenging task, then the model must be encoding relevant information in your predictions. For example, you might think of

training your translation predictions on a question-answering task in the translated language. That is, you use the outputs of your system as inputs to a model for some other task (the question-answering). If your second task can perform as well on your predictions as it can on well-formed data in the translated language, it means that your inputs have the relevant information or patterns for meeting the demands of the task.

The issue with this method is that the second task may not be affected by many of the finer points of translation. For example, if you measured the quality of translation on a query-retrieval task (like pulling up the right webpage for a search query), you would find that a translation that preserves the main topic words of the documents but ignores syntax and grammar might still fit the task well. But this itself doesn't mean that the quality of your translations is accurate or faithful. Therefore, determining the quality of the translation model is just shifted to determining the quality of the task itself, which may or may not be a good standard.

5.3 Bilingual Evaluation Understudy (BLEU)

In 2002, IBM researchers developed the Bilingual Evaluation Understudy (BLEU) that remains, with its many variants to this day, one of the most respected and reliable methods for machine translation.

The BLEU algorithm evaluates the precision score of a candidate machine translation against a reference human translation. The reference human translation is assumed to be a *model* example of a translation, and we use n-gram matches as our metric for how similar a candidate translation is to it. Consider a reference sentence A and candidate translation B:

- A there are many ways to evaluate the quality of a translation, like comparing the number of n-grams between a candidate translation and reference.
- B the quality of a translation is evaluate of n-grams in a reference and with translation.

The BLEU score looks for whether n-grams in the machine translation also appear in the reference translation. Color-coded below are some examples of different size n-grams that are shared between the reference and candidate translation.

- A there are many ways to evaluate the quality of a translation, like comparing the number of n-grams between a candidate translation and reference.

B the quality of a translation is evaluate of n-grams in a reference and with translation.

The BLEU algorithm identifies all such matches of n-grams above, including the unigram matches, and evaluates the strength of the match with the *precision* score. The precision score is the fraction of n-grams in the translation that also appear in the reference.

The algorithm also satisfies two other constraints. For each n-gram size, a gram in the reference translation cannot be matched more than once. For example, the unigram "a" appears twice in B but only once in A. This only counts for one match between the sentences. Additionally, we impose a brevity penalty so that very small sentences that would achieve a 1.0 precision (a "perfect" matching) are not considered good translations. For example, the single word "there" would achieve a 1.0 precision match, but it is obviously not a good match.

Let us see how to actually compute the BLEU score. First let k be the maximum n-gram that we want to evaluate our score on. That is, if $k = 4$, the BLUE score only counts the number of n-grams with length less than or equal to 4, and ignores larger n-grams. Let

$$p_n = \# \text{ matched n-grams} / \# \text{ n-grams in candidate translation}$$

the precision score for the grams of length n . Finally, let $w_n = 1/2^n$ be a geometric weighting for the precision of the n 'th gram. Our brevity penalty is defined as

$$\beta = e^{\min(0, 1 - \frac{\text{len}_{\text{ref}}}{\text{len}_{\text{MT}}})}$$

where len_{ref} is the length of the reference translation and len_{MT} is the length of the machine translation.

The BLEU score is then defined as

$$\text{BLEU} = \beta \prod_{i=1}^k p_n^{w_n}$$

The BLEU score has been reported to correlate well with human judgment of good translations, and so remains a benchmark for all evaluation metrics following it. However, it does have many limitations. It only works well on the corpus level because any zeros in precision scores will zero the entire BLEU score. Additionally, this BLEU score as presented suffers for only comparing a candidate translation against a single reference, which is surely a noisy representation of the relevant n-grams that need to be matched. Variants of BLEU have modified the algorithm to compare the candidate with multiple reference examples. Additionally, BLEU scores may only be

a necessary but not sufficient benchmark to pass for a good machine translation system. Many researchers have optimized BLEU scores until they have begun to approach the same BLEU scores between reference translations, but the true quality remains far below human translations.

6 Dealing with the large output vocabulary

Despite the success of modern NMT systems, they have a hard time dealing with large vocabulary size. Specifically, these Seq2Seq models predict the next word in the sequence by computing a target probabilistic distribution over the entire vocabulary using *softmax*. It turns out that **softmax** can be quite expensive to compute with a large vocabulary and its complexity also scales proportionally to the vocabulary size. We will now examine a number of approaches to address this issue.

6.1 Scaling softmax

A very natural idea is to ask "can we find more efficient ways to compute the target probabilistic distribution?" The answer is Yes! In fact, we've already learned two methods that can reduce the complexity of "softmax", which we'll present a high-level review below (see details in lecture note 1).

1. Noise Contrastive Estimation

The idea of NCE is to approximate "softmax" by randomly sampling K words from negative samples. As a result, we are reducing the computational complexity by a factor of $\frac{|V|}{K}$, where $|V|$ is the vocabulary size. This method has been proven successful in word2vec. A recent work by Zoph et al.⁷ applied this technique to learning LSTM language models and they also introduced a trick by using the same samples per mini-batch to make the training GPU-efficient.

⁷ Zoph et al. 2016, *Simple, Fast Noise-Contrastive Estimation for Large RNN Vocabularies*

2. Hierarchical Softmax

Morin et al.⁸ introduced a binary tree structure to more efficiently compute "softmax". Each probability in the target distribution is calculated by taking a path down the tree which only takes $O(\log|V|)$ steps. Notably, even though Hierarchical Softmax saves computation, it cannot be easily parallelized to run efficiently on GPU. This method is used by Kim et al.⁹ to train character-based language models which will be covered in lecture 13.

⁸ Morin et al. 2005, *Hierarchical Probabilistic Neural Network Language Model*

⁹ Kim et al. 2015, *Character-Aware Neural Language Models*

One limitation for both methods is that they only save computation during training step (when target word is known). At test time,

one still has to compute the probability of all words in the vocabulary in order to make predictions.

6.2 Reducing vocabulary

Instead of optimizing "softmax", one can also try to reduce the effective vocabulary size which will speed up both training and test steps. A naive way of doing this is to simply limit the vocabulary size to a small number and replace words outside the vocabulary with a tag <UNK>. Now, both training and test time can be significantly reduced but this is obviously not ideal because we may generate outputs with lots of <UNK>.

Jean et al.¹⁰ proposed a method to maintain a constant vocabulary size $|V'|$ by partitioning the training data into subsets with τ unique target words, where $\tau = |V'|$. One subset can be found by sequentially scanning the original data set until τ unique target words are detected (Figure 7). And this process is iterated over the entire data set to produce all mini-batch subsets. In practice, we can achieve about 10x saving with $|V| = 500K$ and $|V'| = 30K, 50K$.

This concept is very similar to NCE in that for any given word, the output vocabulary contains the target word and $|V'| - 1$ negative (noise) samples. However, the main difference is that these negative samples are sampled from a biased distribution Q for each subset V' where

$$Q(y_t) = \begin{cases} \frac{1}{|V'|}, & \text{if } y_t \in V' \\ 0, & \text{otherwise} \end{cases}$$

At test time, one can similarly predict target word out of a selected subset, called *candidate list*, of the entire vocabulary. The challenge is that the correct target word is unknown and we have to "guess" what the target word might be. In the paper, the authors proposed to construct a candidate list for each source sentence using K most frequent words (based on unigram probability) and K' likely target words for each source word in the sentence. In Figure 8), an example is shown with $K' = 3$ and the candidate list consists of all the words in purple boxes. In practice, one can choose the following values: $K = 15k, 30k, 50k$ and $K' = 10, 20$.

6.3 Handling unknown words

When NMT systems use the techniques mentioned above to reduce effective vocabulary size, inevitably, certain words will get mapped to <UNK>. For instance, this could happen when the predicted word, usually rare word, is out of the candidate list or when we encounter

¹⁰ Jean et al. 2015, *On Using Very Large Target Vocabulary for Neural Machine Translation*

Sequentially select examples: $|V'| = 5$.

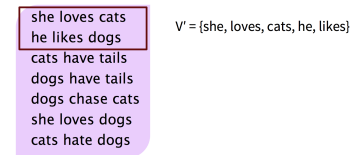


Figure 7: Training data partition

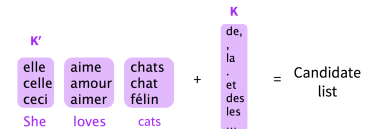


Figure 8: Candidate list

unseen words at test time. We need new mechanisms to address the rare and unknown word problems.

One idea introduced by Gulcehre et al.¹¹ to deal with these problems is to learn to "copy" from source text. The model (Figure 9) applies attention distribution l_t to decide *where* to point in the source text and uses the decoder hidden state S_t to predict a binary variable Z_t which decides *when* to copy from source text. The final prediction is either the word y_t^w chosen by softmax over candidate list, as in previous methods, or y_t^l copied from source text depending on the value of Z_t . They showed that this method improves performance in tasks like machine translation and text summarization.

As one can imagine, there are of course limitations to this method. It is important to point out a comment from Google's NMT paper¹² on this method, "this approach is both unreliable at scale — the attention mechanism is unstable when the network is deep — and copying may not always be the best strategy for rare words — sometimes transliteration is more appropriate".

7 Word and character-based models

As discussed in section 6, "copy" mechanisms are still not sufficient in dealing with rare or unknown words. Another direction to address these problems is to operate at sub-word levels. One trend is to use the same seq2seq architecture but operate on a smaller unit — word segmentation, character-based models. Another trend is to embrace hybrid architectures for words and characters.

7.1 Word segmentation

Sennrich et al.¹³ proposed a method to enable open-vocabulary translation by representing rare and unknown words as a sequence of subword units.

This is achieved by adapting a compression algorithm called **Byte Pair Encoding**. The essential idea is to start with a vocabulary of characters and keep extending the vocabulary with most frequent n-gram pairs in the data set. For instance, in Figure 10, our data set contains 4 words with their frequencies on the left, i.e. "low" appears 5 times. Denote (p, q, f) as a n-gram pair p, q with frequency f. In this figure, we've already selected most frequent n-gram pair (e,s,9) and now we are adding current most frequent n-gram pair (es,t,9). This process is repeated until all n-gram pairs are selected or vocabulary size reaches some threshold.

One can choose to either build separate vocabularies for training and test sets or build one vocabulary jointly. After the vocabulary

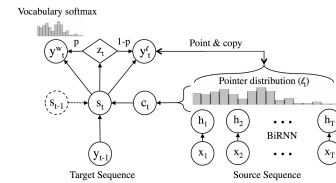


Figure 9: Pointer network Architecture

¹¹ Gulcehre et al. 2016, *Pointing the Unknown Words*

¹² Wu et al. 2016, *Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation*

¹³ Sennrich et al. 2016, *Neural Machine Translation of Rare Words with Subword Units*



Figure 10: Byte Pair Encoding

is built, an NMT system with some seq2seq architecture (the paper used Bahdanau et al. ¹⁴), can be directly trained on these word segments. Notably, this method won top places in WMT 2016.

¹⁴ Bahdanau et al. 2014, "Neural Machine Translation by Jointly Learning to Align and Translate"

7.2 Character-based model

Ling et al. ¹⁵ proposed a character-based model to enable open-vocabulary word representation.

¹⁵ Ling, et al. 2015, "Finding Function in Form: Compositional Character Models for Open Vocabulary Word Representation"

For each word w with m characters, instead of storing a word embedding, this model iterates over all characters $c_1, c_2 \dots c_m$ to look up the character embeddings $e_1, e_2 \dots e_m$. These character embeddings are then fed into a biLSTM to get the final hidden states h_f, h_b for forward and backward directions respectively. The final word embedding is computed by an affine transformation of two hidden states:

$$e_w = W_f H_f + W_b H_b + b$$

There are also a family of CNN character-based models which will be covered in lecture 13.

7.3 Hybrid NMT

Luong et al. ¹⁶ proposed a Hybrid Word-Character model to deal with unknown words and achieve open-vocabulary NMT. The system translates mostly at word-level and consults the character components for rare words. On a high level, the character-level recurrent neural networks compute source word representations and recover unknown target words when needed. The twofold advantage of such a hybrid approach is that it is much faster and easier to train than character-based ones; at the same time, it never produces unknown words as in the case of word-based models.

¹⁶ Luong et al. 2016, *Achieving Open Vocabulary Neural Machine Translation with Hybrid Word-Character Models*

Word-based Translation as a Backbone The core of the hybrid NMT is a deep LSTM encoder-decoder that translates at the word level. We maintain a vocabulary of size $|V|$ per language and use $\langle unk \rangle$ to represent out of vocabulary words.

Source Character-based Representation In regular word-based NMT, a universal embedding for $\langle unk \rangle$ is used to represent all out-of-vocabulary words. This is problematic because it discards valuable information about the source words. Instead, we learn a deep LSTM model over characters of rare words, and use the final hidden state of the LSTM as the representation for the rare word (Figure 11).

Target Character-level Generation General word-based NMT allows generation of $\langle unk \rangle$ in the target output. Instead, the goal here is to create a coherent framework that handles an unlimited output vocabulary. The solution is to have a separate deep LSTM that

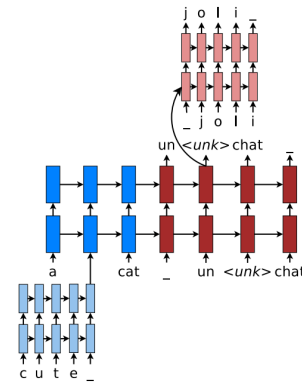


Figure 11: Hybrid NMT

"translates" at the character level given the current word-level state.
Note that the current word context is used to initialize the character-level encoder. The system is trained such that whenever the word-level NMT produces an *<unk>*, the character-level decoder is asked to recover the correct surface form of the unknown target word.

CS224n: Natural Language Processing with Deep Learning¹

Lecture Notes: Part VII

Question Answering²

Winter 2019

¹ Course Instructors: Christopher Manning, Richard Socher

² Authors: Francois Chaubard, Richard Socher

1 Dynamic Memory Networks for Question Answering over Text and Images

The idea of a QA system is to extract information (sometimes passages, or spans of words) directly from documents, conversations, online searches, etc., that will meet user's information needs. Rather than make the user read through an entire document, QA system prefers to give a short and concise answer. Nowadays, a QA system can combine very easily with other NLP systems like chatbots, and some QA systems even go beyond the search of text documents and can extract information from a collection of pictures.

There are many types of questions, and the simplest of them is factoid question answering. It contains questions that look like "The symbol for mercuric oxide is?" "Which NFL team represented the AFC at Super Bowl 50?". There are of course other types such as mathematical questions (" $2+3=?$ "), logical questions that require extensive reasoning (and no background information). However, we can argue that the information-seeking factoid questions are the most common questions in people's daily life.

In fact, most of the NLP problems can be considered as a question-answering problem, the paradigm is simple: we issue a query, and the machine provides a response. By reading through a document, or a set of instructions, an intelligent system should be able to answer a wide variety of questions. We can ask the POS tags of a sentence, we can ask the system to respond in a different language. So naturally, we would like to design a model that can be used for general QA.

In order to achieve this goal, we face two major obstacles. Many NLP tasks use different architectures, such as TreeLSTM (Tai et al., 2015) for sentiment analysis, Memory Network (Weston et al., 2015) for question answering, and Bi-directional LSTM-CRF (Huang et al., 2015) for part-of-speech tagging. The second problem is full multi-task learning tends to be very difficult, and transfer-learning remains to be a major obstacle for current neural network architectures across artificial intelligence domains (computer vision, reinforcement learning, etc.).

We can tackle the first problem with a shared architecture for NLP: Dynamic Memory Network (DMN), an architecture designed for

general QA tasks. QA is difficult, partially because reading a long paragraph is difficult. Even for humans, we are not able to store a long document in your working memory.

1.1 Input Module

Dynamic Memory Network is divided into modules. First we look at input module. The input module takes as input a sequence of T_I words and outputs a sequence of T_C fact representations. If the output is a list of words, we have $T_C = T_I$ and if the output is a list of sentences, we have T_C as the number of sentences and T_I as the number of words in the sentences. We use a simple GRU to read the sentences in, i.e. the hidden state $h_t = \text{GRU}(x_t, h_{t-1})$ where $x_t = L[w_t]$, where L is the embedding matrix and w_t is the word at time t . We further improve it by using Bi-GRU as shown in Figure 2.

1.2 Question Module

We also use a standard GRU to read in the question (using embedding matrix L : $q_t = \text{GRU}(L[w_t^Q], q_{t-1})$), but the output of the question module is an encoded representation of the question.

1.3 Episodic Memory Module

One of the distinctive features of the dynamic memory network is the episodic memory module which runs over the input sequence multiple times, each time paying attention to a different subset of facts from the input.

It accomplishes this using a Bi-GRU that takes input of the sentence-level representation passed in from the input module, and produces an episodic memory representation.

We denote the episodic memory representation as m^i and the episode representation (output by the attention mechanism) as e^i . The episodic memory representation is initialized using $m^0 = q$, and proceeds using the GRU: $m^i = \text{GRU}(e^i, m^{i-1})$. The episode representation is updated using the hidden state outputs from the input module as follows where g is the attention mechanism:

$$h_t^i = g_t^i \text{GRU}(c_t, h_{t-1}^i) + (1 - g_t^i) h_{t-1}^i$$

$$e_i = h_{T_C}^i$$

The attention vector g may be computed in a number of ways, but in the original DMN paper (Kumar et al. 2016), the following formulation was found to work best:

$$g_t^i = G(c_t, m^{i-1}, q)$$

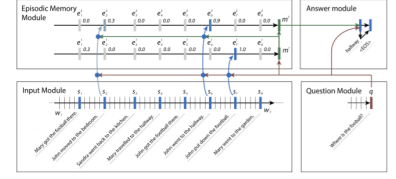


Figure 1: A graphical illustration of the Dynamic Memory Network.

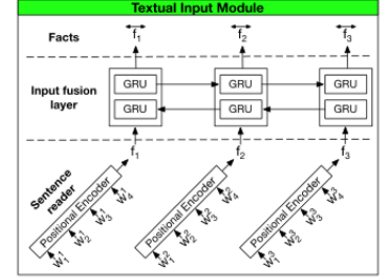


Figure 2: A graphical illustration of the Dynamic Memory Network.

$$G(c, m, q) = \sigma(W^{(2)} \tanh(W^{(1)} z(c, m, q) + b^{(1)}) + b^{(2)})$$

$$z(c, m, q) = [c, m, q, c \circ q, c \circ m, |c - q|, |c - m|, c^T W^{(b)} q, c^T W^{(b)} m]$$

In this way, gates in this module are activated if the sentence is relevant to the question or memory. In the i th pass, if the summary is not sufficient to answer the question, we can repeat sequence over input in the $(i + 1)$ th pass. For example, consider the question "Where is the football?" and input sequences "John kicked the football" and "John was in the field." In this example, *John* and *football* could be linked in one pass and then *John* and *field* could be linked in the second pass, allowing for the network to perform a transitive inference based on the two pieces of information.

1.4 Answer Module

The answer module is a simple GRU decoder that takes in the output of question module, and episodic memory module, and output a word (or in general a computational result). It works as follows:

$$y_t = \text{softmax}(W^{(a)} a_t)$$

$$a_t = \text{GRU}([y_{t-1}, q], a_{t-1})$$

1.5 Experiments

Through the experiments we can see DMN is able to outperform MemNN in babl question answering tasks, and it can outperform other architectures for sentiment analysis and part-of-speech tagging. How many episodes are needed in the episodic memory? The answer is that the harder the task is, the more passes are required. Multiple passes also allows the network to truly comprehend the sentence by paying attention to only relevant parts for the final task, instead of reacting to just the information from the word embedding.

The key idea is to modularize the system, and you can allow different types of input by change the input module. For example, if we replace the input module with a convolutional neural network-based module, then this architecture can handle a task called visual question answering (VQA). It is also able to outperform other models in this task.

1.6 Summary

The zeal to search for a general architecture that would solve all problems has slightly faded since 2015, but the desire to train on one domain and generalize to other domains has increased. To comprehend more advanced modules for question answering, readers can refer to the dynamic coattention network (DCN).